

Nelli, L., Ferguson, H. M. and Matthiopoulos, J. (2019) Achieving explanatory depth and spatial breadth in infectious disease modelling: Integrating active and passive case surveillance. *Statistical Methods in Medical Research*, (doi:[10.1177/0962280219856380](https://doi.org/10.1177/0962280219856380))

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/188671/>

Deposited on: 19 June 2019

Enlighten – Research publications by members of the University of Glasgow
<http://eprints.gla.ac.uk>



Achieving explanatory depth and spatial breadth in infectious disease modelling: integrating active and passive case surveillance.

Journal:	<i>Statistical Methods in Medical Research</i>
Manuscript ID	SMM-18-0591.R2
Manuscript Type:	Original Article
Keywords:	Bayesian modelling, Disease mapping, Imperfect detection, Latent point process, N-mixture models, Spatial epidemiology
Abstract:	<p>Ideally, the data used for robust spatial prediction of disease distribution should be both high-resolution and spatially expansive. However, such in-depth and geographically broad data are rarely available in practice. Instead, researchers usually acquire either detailed epidemiological data with high resolution at a small number of active sampling sites, or more broad-ranging but less precise data from passive case surveillance. We propose a novel inferential framework, capable of simultaneously drawing insights from both passive and active data types. We developed a Bayesian latent point process approach, combining active data collection in a limited set of points, where in-depth covariates are measured, with passive case detection, where error-prone, large-scale disease data are accompanied only by coarse or remotely-sensed covariate layers.</p> <p>Using the example of malaria, we tested our method's efficiency under several hypothetical scenarios of reported incidence in different combinations of imperfect detection and spatial complexity of the environmental variables.</p> <p>We provide a simple solution to a widespread problem in spatial epidemiology, combining latent process modelling and spatially autoregressive modelling. By using active sampling and passive case detection in a complementary way, we achieved the best-of-both-worlds, in effect, a formal calibration of spatially extensive, error-prone data by localised, high-quality data.</p>
<p>Note: The following files were submitted by the author for peer review, but cannot be converted to PDF. You must view these files (e.g. movies) online.</p>	

Achieving explanatory depth and spatial breadth in infectious disease modelling: integrating active and passive case surveillance.

Luca Nelli^{1*}, Heather M. Ferguson¹, Jason Matthiopoulos¹

¹University of Glasgow, Institute of Biodiversity Animal Health and Comparative Medicine, Glasgow G12 8QQ, UK

*corresponding author: luca.nelli@glasgow.ac.uk

Abstract

Ideally, the data used for robust spatial prediction of disease distribution should be both high-resolution and spatially expansive. However, such in-depth and geographically broad data are rarely available in practice. Instead, researchers usually acquire either detailed epidemiological data with high resolution at a small number of active sampling sites, or more broad-ranging but less precise data from passive case surveillance.

We propose a novel inferential framework, capable of simultaneously drawing insights from both passive and active data types. We developed a Bayesian latent point process approach, combining active data collection in a limited set of points, where in-depth covariates are measured, with passive case detection, where error-prone, large-scale disease data are accompanied only by coarse or remotely-sensed covariate layers.

Using the example of malaria, we tested our method's efficiency under several hypothetical scenarios of reported incidence in different combinations of imperfect detection and spatial complexity of the environmental variables.

We provide a simple solution to a widespread problem in spatial epidemiology, combining latent process modelling and spatially autoregressive modelling. By using active sampling and passive case detection in a complementary way, we achieved the best-of-both-worlds, in effect, a formal calibration of spatially extensive, error-prone data by localised, high-quality data.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

25 **Keywords:** Bayesian modelling, disease mapping, imperfect detection, latent point process, N-
26 mixture models, spatial epidemiology.

For Peer Review

28 INTRODUCTION

29 Predictive maps of disease risk, typically obtained by modelling the spatial heterogeneity in disease
30 incidence as a function of underlying covariates, can be crucial for targeting effective control and
31 surveillance ¹⁻⁶. However, reliable prediction at the landscape scale is often hindered by lack of
32 appropriate, high resolution spatial data. Traditionally, incidence data and potential explanatory
33 covariates are collected either systematically – using active sampling by researchers – or
34 opportunistically – from clinical records reported at health facilities. Each of these sampling
35 strategies has its own limitations ⁷. For example, by collecting detailed data for both disease
36 incidence and related covariates, data from active sampling allows models to achieve high
37 explanatory power but not to make large-scale extrapolation and predictions in areas where fine
38 scale covariates are not directly measurable ^{8,9}. On the other hand, passive sampling yields data
39 from a large number of geographically dispersed cases which are more amenable for large scale
40 predictions, but these data often suffer from severe reporting biases ¹⁰⁻¹³ and can be paired with
41 only coarse environmental covariates that have limited explanatory power ⁴. As the drawbacks of
42 one strategy are clearly the strengths of the other, modelling frameworks that consider these two
43 types of data simultaneously and complementarily would strengthen our biological insights and
44 predictive power.

45 Active sampling is typically conducted by research teams that focus on a small number of
46 predetermined locations, with collection of detailed environmental or epidemiological variables
47 including clinical samples ^{9, 14-19}, entomological indicators (for vector-borne disease) ^{17, 18, 20-22}, human
48 demographic and socio-economic factors ^{9, 19, 23, 24} or fine-scale environmental conditions ^{25, 26}. Such
49 data can provide high power for explaining variation in risk across focal sites ⁷, but lack predictive
50 breadth across space because many of the crucial covariates are not available for un-sampled
51 locations ⁹.

1
2
3 52 Clinical records from passive case detection offer the potential of expansive descriptions of spatial
4
5 53 incidence patterns. However, since these incidence data often arise from self-reporting at health
6
7 54 centres, they can be biased by their opportunistic nature. Reporting bias is well acknowledged for
8
9 55 numerous infectious disease systems ²⁷⁻²⁹ and can be expressed as a combined function of distance
10
11 56 from health facilities, the likelihood of asymptomatic cases and sociodemographic factors ^{10-13, 30-35} or
12
13
14 57 more complex measures of travel time ³⁶. Despite this limitation, health centre surveys remain the
15
16 58 primary source of information for disease monitoring. Another drawback of spatial models of
17
18 59 incidence data gathered from passive case detection, relates to the availability of environmental
19
20 60 predictor data. If the locality of the patient is recorded, incidence data can be spatially plotted but
21
22 61 researchers and public health workers are unlikely to be able to directly measure some detailed
23
24 62 explanatory variables at those localities. Therefore, when modelling the incidence data, only large-
25
26 63 scale but coarse layers are customarily considered. While these bring more geographically expansive
27
28 64 information than the highly localised survey data, they generally consist of remotely sensed
29
30 65 covariates and summary records such as bioclimatic, geomorphological, vegetation indexes, human
31
32 66 population density or road networks ^{9, 37-40}, that typically contribute limited explanatory power.
33
34
35
36
37 67 Some studies ^{7, 14-16, 23, 41-44} make use of data from both active and passive case detections together,
38
39 68 but focus on independent analysis and comparison of results from these separate data sources
40
41 69 rather than integrating them. Analysing these two data sources jointly can be viewed as challenging
42
43 70 ¹⁵ because their limitations imply a trade-off between explanatory depth and predictive breadth.
44
45
46 71 However, there is clearly an opportunity to achieve complementarity by analysing them on an
47
48 72 integrated inferential framework. Here, for the first time, we develop a spatial statistical model
49
50 73 combining these two sources of incidence data to harness the maximum amount of information for
51
52 74 explanatory and predictive objectives.
53
54
55
56 75 Our framework takes a novel approach to both the response and the explanatory variables. The dual
57
58 76 nature of the incidence data requires specification of a statistical model that considers two different
59
60

77 aspects of likelihood, one for the localised but precise survey data, and another for the spatially
 78 extensive but imperfect clinical reporting data. We build this part of the approach on two
 79 cornerstones of the statistical literature: the point process model ^{45, 46} and the methodology of point
 80 transects ⁴⁷. Point processes model events (e.g. infection cases) that occur continuously in space
 81 according to an unknown intensity (a spatial surface to be estimated as a function of covariates). We
 82 observe these events as arising from two different point transects, each having its own spatially
 83 heterogeneous observation model. The first type of observation point is the active sampling
 84 location, where cases are detected near-perfectly, but only for that particular set of geographical
 85 coordinates. The second type of observation point is a clinic, where cases of the disease are reported
 86 from a broad geographical region but with probabilities of detection that decay with distance from
 87 the clinic. Regarding the explanatory variables, some environmental variables are easily collectable
 88 for both passive case detection and active sampling points, but more important and powerful
 89 variables may be available only at the latter. By importing ideas from latent process modelling ^{48, 49}
 90 we use the spatially extensive clinical data together with the data-rich survey data to reconstruct
 91 latent covariates that may be hidden from direct or remote observation.

92 To validate the ability of our model to retrieve correct parameter values we require these scenarios
 93 to be accompanied by known intensity surfaces for both incidence and latent explanatory variables.
 94 These requirements cannot be satisfied by real data sets, so here we have acquired our scenarios via
 95 realistic simulation, motivating our examples from a real system of a vector-borne disease. To
 96 illustrate the generality of our approach, we have hypothesized multiple contrasting scenarios of
 97 reporting bias and spatial distribution of the latent process underlying disease incidence.

98 We chose malaria in West Africa as an ideal example of an important environmentally-dependent
 99 infectious disease ^{50, 51}, for which human exposure and infection risk is highly spatially
 100 heterogeneous and dependent on crucial environmental variables that influence interactions
 101 between people, mosquitoes and parasites ^{40, 52}. Control measures such as long-lasting insecticide

treated nets (LLINs) have been crucial for impeding contact between mosquitoes and people, and have led to substantial declines in malaria prevalence across Africa in the last decade^{50, 51}. However, the success of such an approach may be undermined by development of insecticide resistance in mosquitoes, particularly in West Africa where rates are among the highest in the world⁵³⁻⁵⁵. Copious and widespread data on reported cases are often available from clinics (see for example www.malariasurveys.org or www.dhsprogram.com), but detailed information on mosquito vector ecology and insecticide resistance is only available for a limited number of sites (e.g.^{54, 56-58}). These challenges exist for many other vector-borne diseases whose transmission is dependent on an ecological reservoir and rely on insecticide use for control, such as for example dengue, Zika and chikungunya viruses⁵⁹, Lyme and other tick-borne disease⁶⁰, schistosomiasis⁶¹, Rift Valley Fever⁶², human African trypanosomiasis⁶³ or West Nile Virus⁶⁴.

METHODS

Modelling approaches

For a given area of interest subdivided into a regular grid, we consider as our sampling unit the grid cell $i \in \{1, \dots, K\}$. We first assume an underlying stochastic process f that generates numbers of cases N_i according to an underlying, spatially heterogeneous rate λ_i . We also assume an observation process g that allows a subset of the N_i cases to be reported at different sampling stations. We distinguish between two types of sampling stations: S is the number of active sampling points (about which we are assuming a perfect and exclusive detection but at a small distance, i.e. within the cell that contains them). We denote by J the number of clinics (about which we are assuming an imperfect but long-ranging detection). The observation process g is therefore generating the vector of incidence data reported in each i^{th} cell at different stations $\mathbf{I}_i = \{I_{1,i}, \dots, I_{S,i}, I_{(S+1),i}, \dots, I_{(S+J),i}, \mathbf{U}_i\}$, given the vector of probabilities $\mathbf{P}_i = \{P_{1,i}, \dots, P_{S,i}, P_{(S+1),i}, \dots, P_{(S+J),i}, \mathbf{Q}_i\}$. \mathbf{U}_i represents the number of completely unreported cases in each i^{th} cell (which is a missing value in the data), given the probability \mathbf{Q}_i of not reporting.

The general likelihood function of our models can be expressed as follows:

$$L = \prod_{i=1}^K f(N_i | \lambda_i) g(I_i | \mathbf{P}_i N_i) \quad [1]$$

We built our approach incrementally, developing three distinct modelling approaches with an increasing level of complexity to allow comparison between the routes that might have traditionally been followed to analyse data arising from active sampling (model 1) and passive case detection (model 2) with our new proposed route (model 3), which reconstructs the latent processes and estimates the emergent patterns of disease incidence with increased precision and accuracy.

Model 1 – active sampling data only

Here, we consider data that would be collected from active sampling at just a limited number S of active survey sites. To analyse the relationship between disease incidence and detailed measures of covariates at a set of predetermined survey points, model 1 takes the form of a Poisson Generalised Linear Model without any spatially explicit component.

Although this is a straightforward model to fit using likelihood-based libraries in all statistical platforms, we fitted it using Bayesian methods⁶⁵ for consistency in the comparison with models 2 and 3 that follow. The response variable is the number of observed diseases cases N_i at the location of the i^{th} survey. We assume here (for simplicity, but with no loss of generality) that all the cases at the survey location are recorded (hence, a local detection probability of 1 for each case), although we acknowledge that with conventional diagnostic tests some percentage of cases can be missed⁶⁶. If data are available on diagnostic sensitivity and specificity, our method can be readily extended by incorporating false negatives or positives.

The model takes the form

$$N_i \sim \text{Poisson}(\lambda_i) \quad [2]$$

1
2
3
4
5 147 where the rate (λ_i) of disease incidence is
6
7
8
9
10
11
12
13 148 The linear predictor on the right-hand-side of this expression comprises a set of n coefficients β and
14
15 149 n explanatory variables X measured at the i^{th} survey location.
16
17
18 150 Equations [2] and [3] can be generalised to take better account of specific features of the data. For
19
20 151 example, it may be relevant to use overdispersed forms of the likelihood (relaxing the Poisson
21
22 152 assumption) or more complicated functional forms of the linear predictor, involving polynomials,
23
24 153 interactions or splines.
25
26
27 154 *Model 2 – passive case detection only*
28
29
30 155 Here, we considered only data coming from passive case detection. This model maintained the basic
31
32 156 structure of model 1, i.e. it is a Bayesian Poisson regression, with reported disease cases at human
33
34 157 dwellings or communities surrounding the health centres as the response variable and the set of
35
36
37 158 environmental variables as predictors. Under our scenarios, we assumed that one of the key
38
39 159 predictor variables (insecticide resistance *IR*, see model validation) could only be measured
40
41 160 experimentally in active sampling sites, therefore we couldn't include it in eq. [3].
42
43
44 161 We introduced the estimation of bias in reporting disease cases given by the distance from the
45
46 162 health centres, borrowing concepts from distance sampling theory ⁴⁷, a group of methods, widely
47
48
49 163 used to estimate the absolute abundance or spatial density of animal or plant populations. The key
50
51 164 underlying concept is the estimation of a detection function ($P(d)$), which represents the decay in
52
53 165 the probability of detecting an object with increasing distance (d) from the observer. Given the
54
55 166 detection function and encounter rate, the absolute density of a population can be modelled at a
56
57
58 167 given point, assuming perfect detection at the location of the observer $P(0) = 1$. In our application,
59
60 168 this has the interpretation that if a case arises in the immediate vicinity of the clinic ($d \cong 0$), then it is

certain to be reported. A plausible, but flexible decay function is fitted to paired data of detections and distances. For example, detection of a malaria case from the i^{th} location at the j^{th} clinic, and can be modelled as a half-normal function of distance from the health centre $d_{i,j}$, by the following ⁴⁷:

$$p(d_{i,j}) = \exp\left(-\frac{d_{i,j}^2}{2\sigma^2}\right) \quad [4]$$

where σ is the shape parameter of the half-normal function (regulating how quickly the detection probability drops with distance). The distance d can be Euclidean, or a more complicated function of accessibility (e.g. affected by proximity between points along a given road network).

Any given case may be reported to any one of the available clinics, but clinics nearby are more likely to receive the report. The probability of any one case being reported to any one clinic (accounting for other clinics) can be modelled in terms of the distances of all the clinics from the point of occurrence of the case, as follows

$$P_{i,j} = \frac{p(d_{i,j})}{\sum_{j=1}^J p(d_{i,j}) + Q_i} \quad [5]$$

The denominator here represents all possible outcomes, i.e. the probabilities that the case is reported to any one of J centres, and the probability P_{Q_i} that the case goes completely unreported:

$$Q_i = \prod_{j=1}^J [1 - p(d_{i,j})] \quad [6]$$

Note that $P_{i,j}$ is the standardised form of $p(d_{i,j})$. In fact, $p(d_{i,j})$ is the probability of a case being reported at a given clinic (considered in isolation), purely as a function of distance, whereas $P_{i,j}$ is the probability of reporting at a clinic, accounting for the effects of other clinics that are “competing” for the same reports and including Q_i , that is the probability of a case not being reported at all.

The likelihood of a data set comprising clinic reports may then be written as a multinomial process. In particular, for a given number of actual cases N_i (see eq. [2]), the likelihood of reported disease cases I_i in the i^{th} cell for the J clinics in the dataset is determined by the detection probabilities P_i that are function of distances between the i^{th} location and the clinics, by

$$I_i \sim \text{Multinomial}(N_i, P_i) \quad [7]$$

where $P_i = \{P_{1,i}, \dots, P_{J,i}, Q_i\}$.

Fitting model 2 to the data yielded estimates of the shape parameter of the detection function (eq. [4]) and parameters of eq. [3]. Although it had no spatially explicit component, we used model 2 to generate a reconstruction of the patterns of incidence across space, based on the coarse-level environmental covariates. Hence this model did not benefit from the fine-resolution covariates that could only be measured by detailed experimental methods at survey points.

Model 3 – active and passive data combined

The process and observation model for this joint approach to data took the form of eqs. [2] and [7] respectively. However, just like in model 1, eq. [3] used the full set of predictors, including the partly-latent variable (i.e. insecticide resistance, available only for active sampling points but not for regions of passive case detection data collection and the rest of space). Our model for the latent variable IR postulated a spatial autocorrelation structure⁶⁷, implying that even though we may not know the values of the latent variable at two points in space, we can express a relationship about their expected degree of similarity. Any pair of K cells in our grid, say $i \in \{1, \dots, K\}$ and $k \in \{1, \dots, K\}$, were assumed to have a covariance, specified as a decreasing function of their distance

$$cov_{i,k} = \exp(-\rho d_{i,k}) \quad [8]$$

With $\rho \geq 0$. Again, this is one of many possible structures and our overall approach is not constrained to this functional form.

The distribution of the latent variable $IR = \{IR_1, \dots, IR_K\}$ in all the K cells, was therefore modelled as a Gaussian field from an m -dimensional multivariate normal distribution, where each of the dimensions represented the probability density of a cell in space.

$$IR_i \sim MVN(\mu, \Sigma) \quad [9]$$

Here, the mean vector μ has length K (the total number of cells in geographical space), and Σ is a $K \times K$ spatial covariance matrix⁶⁸ with values of 1 on the diagonal and values $cov_{i,k}$ for the i row and k column from eq. [8].

Model 3, hence, is fitted exactly as model 2 according to eq. [7], but the linear predictor function (eq. [3]), included all the covariates, unlike model 2, which was missing the covariate of IR. In particular, IR observations were used where available (at active sampling points), assuming that they were realisations from eq. [9].

Model validation

We used simulated data on malaria incidence and insecticide resistance within the primary mosquito vectors to validate our models. Our specific validation aims were to 1) evaluate the match between the posterior distribution of the coefficients and the simulation process that generated the data; 2) estimate bias in reporting the clinical data as a function of distance between the location of a clinic and the village where the patient resides; 3) recreate the missing covariate of insecticide resistance \hat{IR} and to reconstruct the true incidence \hat{N} .

Our simulation borrowed its setting from a study currently ongoing in Southwest Burkina Faso (MiRA – Malaria in Insecticide Resistant Africa, Wellcome Trust 200222/Z/15/Z). The study covers an area of approx. 6000 km² in the health district of Banfora in south-western Burkina Faso, comprising primarily West Sudanian savannah which experiences a rainy season from May to October with little

1
2
3 228 rain in other months. Malaria transmission is stable throughout the year but peaks from May to
4
5 229 November. The major vectors are *Anopheles gambiae* and *An. funestus*. Like many other areas of
6
7 230 Africa, the primary malaria control strategy is long lasting insecticidal nets (LLINs) that are
8
9 231 distributed at high coverage across the country (Burkina Faso National Malaria Control Program,
10
11 232 *unpublished data*). In contrast to some areas of Africa, recent LLIN distribution campaigns have had
12
13 233 little impact on malaria prevalence and it is hypothesized that this may be due to high levels of
14
15 234 insecticide resistance in local vector populations ⁶⁹, which are amongst the highest on record.
16
17 235 Resistance to pyrethroid insecticides is widespread. Mortality after exposure (defined by the World
18
19 236 Health Organization (WHO) as the response to the stipulated discriminating dose of permethrin)
20
21 237 ranges from 5-20% ²⁰. For the purposes of data simulation we assume that active sampling of malaria
22
23 238 infections and insecticide resistance levels is carried out in 12 villages, and that passive case data is
24
25 239 available from patients reporting to from 8 health centres distributed throughout the study area.
26
27 240 This number and distribution of passive and active sampling site was selected to represent the
28
29 241 distribution of health facilities and likely maximum amount of active survey data available.
30
31 242 For the simulation, we considered a square grid with a 1km² resolution covering the study area. We
32
33 243 generated a dataset with reported incidence in each cell of the grid under a binomial *N*-mixture
34
35 244 model ^{70, 71} by combining two different processes: a state model, i.e. the biological process that
36
37 245 generates malaria infection cases, and an observation model, i.e. the process that affects the
38
39 246 probability that infection cases are reported to a health centre.
40
41
42 247 To simulate the biological process, we considered the average altitude in the cell, average yearly
43
44 248 temperature (TEMP), annual rainfall (RAIN), human density (HUM), normalised difference vegetation
45
46 249 index (NDVI) and insecticide resistance (IR) in mosquitoes as potential predictors ^{9, 38, 72-77}.
47
48
49 250 Temperature and rainfall were derived from the WorldClim database (www.worldclim.org). NDVI
50
51 251 values were obtained using the package *MODISsp* for R ⁷⁸. To create the layer of human density, we
52
53 252 used a kernel density estimation ⁷⁹ using GPS points of the villages (307) in the study area and the
54
55
56
57
58
59
60

population census in each village (1755 ± 1804 mean \pm dev. std., Institut national de la statistique et de la démographie, *unpublished data*) as weight field. Kernel bandwidth was chosen so as to minimize the least-squares cross validation score (h_{iscv})⁸⁰.

Insecticide resistance reporting has improved over time, and global maps of insecticide resistance at coarse resolutions are now becoming available⁷⁷. However, little is known about its spatial distribution at local scale⁸¹. Therefore, to explore out model's ability to retrieve latent variables of differing spatial complexity, insecticide resistance was simulated by hypothesizing 3 different scenarios of increasing spatial autocorrelation, with parameter ρ of eq. [8] set respectively to $\rho_1 = 3.0$, $\rho_2 = 0.7$ and $\rho_3 = 0.3$ (Fig. 1, *IR1*, *IR2*, *IR3*).

The number of malaria cases, or true incidence, in each cell (N_i) was assumed to have a positive relationship with temperature^{9, 74}, rainfall^{9, 73, 74}, human density⁹, NDVI^{9, 37, 38, 72} and insecticide resistance⁷⁶, and was simulated from eq. [2] using the linear predictor

$$\log(\lambda_i) = \beta_0 + \beta_{HUM}HUM_i + \beta_{NDVI}NDVI_i + \beta_{RAIN}RAIN_i + \beta_{TEMP}TEMP_i + \beta_{IR}IR_i$$

We set the equation's coefficients to the values $\beta_0 = 2.90$, $\beta_{HUM} = 0.50$, $\beta_{NDVI} = 0.30$, $\beta_{RAIN} = 0.20$, $\beta_{TEMP} = 0.25$, $\beta_{IR} = 0.50$. Having 3 distinct scenarios of insecticide resistance *IR1*, *IR2* and *IR3* we obtained 3 scenarios of malaria infection cases $N1_i$, $N2_i$ and $N3_i$.

For the observation process, we accounted for simulated bias in reporting cases in each cell of the grid, by considering a probability of reporting as a function of the distance between a given cell and each health centre. We set the detection probabilities in each cell $P(i,j)$ in accordance with eq. [4] with $p(d_{i,j})$ being the Euclidean distance between the centroid of the i^{th} cell of the grid and each j^{th} health centre. We employed 3 different shapes of the detection function, using different values of the shape parameter $\sigma_A = 10$, $\sigma_B = 15$, $\sigma_C = 20$ (Fig.1, P_A , P_B and P_C). Probability of reporting at active sampling stations was deliberately set at 1, to ensure that all the infection cases occurring at the sampling stations were recorded.

1
2
3 277 By combining the three scenarios of disease incidence given by the biological process with the three
4
5 278 scenarios of detection function, we generated nine different scenarios of reported Incidence for
6
7 279 each cell (I_{ij}), under a multinomial process given by [7]. For each combination scenario the response
8
9 280 data comprised the number of reported cases per cell (Fig. 1, $I1A$ to $I3C$).
10
11
12
13 281 Preliminary manipulation of environmental layers was done using the software QGIS ⁸², the
14
15 282 simulations were conducted in the statistical environment R ⁸³, and Bayesian model fitting to the
16
17 283 simulated data was carried out using the program JAGS ⁸⁴, interfaced with R via the package *rjags* ⁸⁵.
18
19
20 284 We analysed the simulated incidence data, using each of the three models described above. We
21
22 285 used Markov Chain Monte Carlo (MCMC) algorithms (code provided in Appendix S1) to fit each of
23
24 286 the models to the combination of environmental and incidence data. Relatively non-informative
25
26 287 priors were chosen for all process and observation parameters and for the cells of the map relating
27
28 288 to the latent variable. To make this a conservative test of the methodology, we employed priors
29
30 289 wide variances. For the coefficients of the environmental covariates we chose diffuse normal priors
31
32 290 centred at zero, corresponding to a null hypothesis of no-effect for each covariate. For the distance
33
34 291 decay parameter σ of the detection function, we adopted a uniform prior with limits 0-1000 ⁷¹. For
35
36 292 parameter ρ of the covariance matrix describing spatial autocorrelation in the latent covariate, we
37
38 293 used a gamma prior (shape = 0.1, rate = 0.1). To achieve convergence, model 1 and 2 were run for 3
39
40 294 $\times 10^4$, whereas model 3 was run for 1.2×10^6 iterations.
41
42
43
44
45 295 Means of posterior distributions with corresponding credible intervals were obtained for each model
46
47 296 coefficient $\hat{\beta}_k$ as well as the shape parameters of the detection function $\hat{\sigma}$, (only relevant for models
48
49 297 2 and 3). For each model and each simulated scenario, we generated spatial predictions of the
50
51 298 expected true incidence \hat{N} and the latent covariate of insecticide resistance \hat{IR} . The accuracy of each
52
53 299 parameter in the complete set $\theta = (k, \sigma)$ was examined by calculating its relative bias from the true
54
55 300 underlying value, as
56
57
58
59
60

$$RB_{\theta} = \frac{\hat{\theta} - \theta}{|\theta|} \quad [9]$$

and by plotting the simulated vs reconstructed malaria incidence (for models 2 and 3) and between the simulated and reconstructed insecticide resistance (for model 3)

RESULTS

The full results with posterior summaries for all model parameters are reported in the supplementary material (S2). Plots showing the relationship between the simulated and reconstructed malaria incidence and between the simulated and reconstructed insecticide resistance are also presented in supplementary material (S3). Here, we present an overview of these detailed results, by reporting on the values of relative bias $|RB|$ for each explanatory variable, in each model, under the nine different scenarios of reported malaria incidence (Fig. 2).

Model 1 considered only the active sampling points, hence the single column under model 1 in Fig. 2 does not include extended results pertaining to the clinic detection function (see supplementary material S2.1 for full results). Under model 1, the simulated malaria incidence was affected only by the environmental covariates (that were common to all scenarios) including insecticide resistance. Overall, the results from model 1 showed an average $|RB| = 0.11$ (std. dev. = 0.08). This was a persistent finding across all three simulated patterns for the latent variable (IR), with low values of relative bias arising regardless of the degree of spatial autocorrelation of the simulated insecticide resistance layer.

Model 2, which considered only data from passive case detection, was less able to capture the underlying effects of predictors on the reported malaria incidence (see supplementary material S2.2 for full results). The posterior means of all coefficients showed an overall average $|RB| = 0.89$ (std. dev. = 1.52). A pattern of increasing bias emerged in particular when considering scenarios of increasing spatial autocorrelation in the latent variable of insecticide resistance (Fig. 2). Since model 2 only included the passive detection cases, the latent variable was completely missing from the list

1
2
3 324 of covariates. In scenarios *I1A*, *I1B* and *I1C*, given by the same *IR1*, (low spatial autocorrelation),
4
5 325 the average $|RB|$ was 0.87 (std. dev. = 1.53). Scenarios that assumed an intermediate level of spatial
6
7 326 autocorrelation in insecticide resistance (latent variable *IR2*) generated an average $|RB|$ of 0.89
8
9
10 327 (std. dev. = 1.54) whereas models assuming the most spatially autocorrelated distribution of
11
12 328 insecticide resistance (*IR3*) generated an average $|RB|$ of 0.91 (std. dev. = 1.50). Contrary to the
13
14 329 coefficients of the process model, posteriors pertaining to the observation model were not sensitive
15
16 330 to the different shapes of the detection function (cases P_A , P_B or P_C). Posteriors for the parameter $\hat{\sigma}$
17
18 331 of the detection function were highly accurate, with absolute values of relative biases ranging from
19
20 332 0.06 to 0.09 (Fig. 2). This model was able to partly reconstruct disease incidence, but not in areas
21
22
23 333 with relatively higher levels of insecticide resistance (Fig. 3a).
24
25
26 334 Model 3 gave the best results in terms of estimating coefficients with low relative biases (see
27
28 335 supplementary material S2.3 for full results). Of particular note is the fact that the parameter for the
29
30 336 latent insecticide resistance variable $RB_{\hat{\beta}_{IR}}$ showed a low $|RB|$ varying between 0.02 and 0.08.
31
32
33 337 Overall, the average $|RB|$ across all variables was 0.07 (std. dev. = 0.07). As with model 1, but in
34
35 338 contrast to model 2, the magnitude of bias in estimated parameters was unrelated to the degree of
36
37 339 spatial autocorrelation assumed in the latent variable. Similar to model 2, the parameter associated
38
39 340 with the case detection function ($\hat{\sigma}$) was estimated with good accuracy, but model 3 was more
40
41 341 accurate in mapping case distribution (Fig. 3a, see also comparison of plots in supplementary
42
43 342 materials S3.1 vs S3.2). Additionally, the latent distribution of the layer of insecticide resistance was
44
45
46 343 accurately reconstructed using model 3 (Fig. 3b, and supplementary material S3.3).
47
48

49 344 **DISCUSSION**

50
51
52 345 By analysing a wide range of plausible, simulated data sets of disease incidence and environmental
53
54 346 variables arising from active sampling and passive case detection, we uncovered some of the
55
56 347 disadvantages of analysing these two data types in isolation. Additionally, we propose a novel
57
58
59 348 modelling framework aimed at achieving complementarity between the two. We found that such an
60

integrated, spatially-explicit model, which acknowledges both active sampling and passive case detection, leads to great improvements in precision and accuracy but also enables the reconstruction of maps for the hidden variable across unsurveyed space.

As expected, when modelling data arising only from active sampling, we achieved high explanatory power and relatively low bias, because the model had access to measurements of all the covariates underlying disease incidence. The model considering only data coming from passive case detection allowed us to estimate the map of malaria incidence with high accuracy. However, posterior distributions for most parameters were biased which was likely due to missing data for the important variable of insecticide resistance. This condition reflects a common situation in epidemiological studies, where passive case detection at health centres can provide a large amount of long-term data with relatively moderate effort. Our simultaneous estimation of detection functions as part of model inferences shows how to take account of imperfect reporting which is an integral characteristic of such opportunistic data^{12, 27-29, 35}.

With our proposed 3rd model, we achieved a good synergy between depth and breadth in inference by combining the strengths of the first two models, and allowing them to compensate for each other's limitations. In contrast to model using only passive case detection, our hybrid modelling framework allowed us to investigate the effect of all the variables (including the latent one), and to produce accurate predictive maps of the disease incidence and latent variable which were not possible with the model considering only active sampling. An important achievement of our proposed model was the capability to deal with a latent variable, regardless of its level of spatial autocorrelation. Thus, even in the absence of assumptions or any preliminary information on the spatial structure of the latent variable (e.g. whether it is akin to uncorrelated "background noise" or has a highly geography-dependent distribution) this model framework has potential to reconstruct it.

1
2
3 373 Our incremental approach showed that the gains in the accuracy of the results, moving from model
4
5 374 1 to model 3, were a direct result of increases in the spatial complexity used by the analytical
6
7 375 approaches. Model 1 had no explicit spatial component. Model 2 was used to generate predictions
8
9 376 in space but it didn't explicitly consider spatial structure in its formulation. Model 3, by including the
10
11 377 spatial autocorrelation structure in the partly latent variable, led to the best results.
12
13
14
15 378 Our approach to latent variables, readily generalises to processes other than insecticide resistance.
16
17 379 We chose this particular example of a latent variable, because *IR* has potential to impact the
18
19 380 transmission and control of a wide range of vector-borne diseases, including malaria, but is typically
20
21 381 labour-intensive, time-consuming and expensive to measure ⁸⁶. Although WHO guidelines classify
22
23 382 insecticide resistance in a binary way ⁸⁶, the raw data from Tube test bioassays measure the %
24
25 383 survival of cohorts of similarly aged females after a given time period of exposure to insecticide
26
27 384 treated surfaces. Therefore, to greatly increase the inferential value acquired from such data, we
28
29 385 treated insecticide resistance as a continuous variable ranging from 0 to 1. Our approach can be
30
31 386 easily extended to more specific measures of insecticide resistance, such as metabolic, cuticular and
32
33 387 behavioural resistance ^{53, 87}, or to other types of predictor data that can be collected in the field
34
35 388 through active sampling but are not easily obtainable via passive case detection, such as vector
36
37 389 abundance and density.
38
39
40
41
42 390 When simulating and modelling the latent variable, we made an assumption of stationarity (the
43
44 391 autocorrelation function didn't change in space or in time) and monotonicity (the autocorrelation
45
46 392 always decreased with distance). These two assumptions can be plausibly relaxed extending our
47
48 393 autocorrelation function. For example, non-stationary formulations could be achieved by expressing
49
50 394 the rate of autocorrelation decay (ρ_0) as a function of latitude and longitude or time. Alternatively,
51
52 395 ρ_0 could be expressed as log-linear combination of environmental covariates. Non-monotonic
53
54 396 formulations of the autocorrelation function could be produced for cases where periodic patterns
55
56
57
58
59
60

exist in space, but we currently see very little justification for such formulations based on biological first principles.

The ability to account for reporting bias of our response variable, makes our approach easily applicable to other scenarios where an imperfect detection needs to be considered, such as citizen science data⁸⁸ or mobile phone surveillance tools⁸⁹. When modelling the detection function, we made similar assumptions (stationarity and monotonicity) to those of the autocorrelation function for the latent variable and we hypothesized the observation process was only affected by distance from health centres^{10-12, 30, 32, 34}. In several real-world scenarios, additional covariates of reporting probability may be involved, such as age and sex of the patient and socioeconomic factors^{31, 33, 35}. Borrowing fundamental concepts from Distance sampling⁴⁷, we assumed that at zero distance the probability of reporting the disease was 100%, however asymptomatic disease in apparently healthy people is common^{66, 90}, and would not be observed in clinical data. Thus, incomplete detection at zero distance (based on additional calibration data on the frequency of asymptomatic cases) must be considered⁹¹. Finally, human mobility is unlikely to be strictly related to Euclidean distance (a third implicit assumption of our detection function), so it may be preferable to use the distance according to road network⁶, when applying this model to real data. Global digital layers quantifying underlying “landscape resistance”, describe the travel time between any two points on the globe, based on data such as road density, terrain morphology and an political borders³⁶, could be easily included in a spatially explicit epidemiological model such as ours⁹². For all of these reasons, we suggest that preliminary analysis using pilot data and focussing only on modelling the detection probability should be carried out before integrating it into the final model.

Our likelihood could be deployed using either a Bayesian or a frequentist setting. It is likely that in real life, most epidemiological data sets will be accompanied by sufficient expert opinion to lead to influential priors, hence we have illustrated using a Bayesian approach. However, we did not assume the existence of expert opinion here, because we were seeking to construct a conservative test of

1
2
3 422 our methods. The models presented here (in particular our model 3, using both data types) require a
4
5 423 high computational effort (see supplementary material for details). Notwithstanding their
6
7 424 theoretical simplicity, the need to take spatial structure into account with a large dataset slows
8
9 425 down the Bayesian MCMC inference. Other model fitting approaches such as the Integrated Nested
10
11 426 Laplace approximation (INLA)⁹³, may prove capable of providing similarly accurate results but with
12
13 427 faster processing ⁹⁴.
14
15
16
17 428 In quantitative ecology, data simulation, by generating random realisations from stochastic
18
19 429 processes described by a series of distributional statements, is exceedingly useful ⁷¹. Although
20
21 430 simulated studies are not guaranteed to be the same as a real epidemiological system, they allow
22
23 431 objective validation of proposed frameworks on a wide range of plausible scenarios, easily adaptable
24
25 432 to other epidemiological studies. Although our simulation was borrowing its settings from a study
26
27 433 specifically looking at malaria, we demonstrated its applicability on a broad range of contrasting
28
29 434 scenarios. Therefore we believe that such a framework can successfully work under different
30
31 435 epidemiological systems, where a combination of large-scale but opportunistic data are collected at
32
33 436 the same time as conducting a small number of localised scientific surveys.
34
35
36
37 437 The strength of our proposed analytical approach lies in its ability to use distinct solutions, such as
38
39 438 latent process modelling and spatially autoregressive modelling, in a fully integrated framework. In
40
41 439 particular, we demonstrated how active sampling and passive case detection, that have so far been
42
43 440 considered independently in the context of spatial epidemiology, can be used simultaneously and
44
45 441 complementarily in a package where the strength of one compensates for the drawback of the other.
46
47 442 Our method shows promise for complex spatial epidemiology studies, by allowing different parts of
48
49 443 the model to glean information from different types of data. Such egalitarian and complementary
50
51 444 use of two, or more data types, can be extended to make use of digital or hard copy primary care
52
53 445 records, irrespective of the sophistication of the health provision systems, the density of the human
54
55 446 population, or the nature of the disease.
56
57
58
59
60

447

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Acknowledgements

We would like to thank Laurie Baker, Julie Barker, Fergus Chadwick, Heather McDevitt and Hilary Ranson for their useful comments on an earlier draft of the manuscript. We also thank the two anonymous reviewers who provided essential suggestions to improve the presentation of the statistical elements.

Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

This project has received funding from the Wellcome Trust under grant agreement no [200222/Z/15/Z] MiRA.

For Peer Review

REFERENCES

1. Pfeiffer D, Robinson TP, Stevenson M, Stevens KB, Rogers DJ and Clements AC. *Spatial analysis in epidemiology*. Oxford University Press New York, 2008.
2. Reisen WK. Landscape epidemiology of vector-borne diseases. *Annual review of entomology*. 2010; 55: 461-83.
3. Woolhouse M. How to make predictions about future infectious disease risks. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*. 2011; 366: 2045-54.
4. Hartemink N, Vanwambeke SO, Purse BV, Gilbert M and Van Dyck H. Towards a resource-based habitat approach for spatial modelling of vector-borne disease risks. *Biological Reviews*. 2015; 90: 1151-62.
5. Lawson AB, Banerjee S, Haining RP and Ugarte MD. *Handbook of spatial epidemiology*. CRC Press, 2016.
6. Kirby RS, Delmelle E and Eberth JM. Advances in spatial epidemiology and geographic information systems. *Annals of Epidemiology*. 2017; 27: 1-9.
7. Bakuza JS, Denwood MJ, Nkwengulila G and Mable BK. Estimating the prevalence and intensity of *Schistosoma mansoni* infection among rural communities in Western Tanzania: The influence of sampling strategy and statistical approach. *PLOS Neglected Tropical Diseases*. 2017; 11: e0005937.
8. Eisen L and Eisen RJ. Using Geographic Information Systems and Decision Support Systems for the Prediction, Prevention, and Control of Vector-Borne Diseases. *Annual Review of Entomology*. 2011; 56: 41-61.
9. Samadoulougou S, Maheu-Giroux M, Kirakoya-Samadoulougou F, De Keukeleire M, Castro MC and Robert A. Multilevel and geo-statistical modeling of malaria risk in children of Burkina Faso. *Parasites & vectors*. 2014; 7: 350.
10. Müller I, Smith T, Mellor S, Rare L and Genton B. The effect of distance from home on attendance at a small rural health centre in Papua New Guinea. *International journal of epidemiology*. 1998; 27: 878-84.
11. Feikin DR, Nguyen LM, Adazu K, et al. The impact of distance of residence from a peripheral health facility on pediatric health utilisation in rural western Kenya. *Tropical Medicine & International Health*. 2009; 14: 54-61.
12. Biswas RK and Kabir E. Influence of distance between residence and health facilities on non-communicable diseases: An assessment over hypertension and diabetes in Bangladesh. *PLoS ONE*. 2017; 12: e0177027.
13. Minuzzi-Souza TTC, Nitz N, Cuba CAC, et al. Surveillance of vector-borne pathogens under imperfect detection: lessons from Chagas disease risk (mis)measurement. *Scientific Reports*. 2018; 8: 151.
14. Tiono AB, Kangoye DT, Rehman AM, et al. Malaria incidence in children in South-West Burkina Faso: comparison of active and passive case detection methods. *PLoS One*. 2014; 9: e86936.
15. Bhoomiboonchoo P, Nisalak A, Chansatiporn N, et al. Sequential dengue virus infections detected in active and passive surveillance programs in Thailand, 1994–2010. *BMC Public Health*. 2015; 15: 250.
16. Sarti E, L’Azou M, Mercado M, et al. A comparative study on active and passive epidemiological surveillance for dengue in five countries of Latin America. *International Journal of Infectious Diseases*. 2016; 44: 44-9.
17. Pham Thi KL, Briant L, Gavotte L, et al. Incidence of dengue and chikungunya viruses in mosquitoes and human patients in border provinces of Vietnam. *Parasites & Vectors*. 2017; 10: 556.
18. Fauver JR, Weger-Lucarelli J, Fakoli LS, III, et al. Xenosurveillance reflects traditional sampling techniques for the identification of human pathogens: A comparative study in West Africa. *PLOS Neglected Tropical Diseases*. 2018; 12: e0006348.

19. Mai VQ, Mai TTX, Tam NLM, Nghia LT, Komada K and Murakami H. Prevalence and Risk Factors of Dengue Infection in Khanh Hoa Province, Viet Nam: A Stratified Cluster Sampling Survey. *Journal of Epidemiology*. 2018; advpub.
20. Bagi J, Grisales N, Corkill R, et al. When a discriminating dose assay is not enough: measuring the intensity of insecticide resistance in malaria vectors. *Malaria Journal*. 2015; 14.
21. Krajacich BJ, Slade JR, Mulligan RF, et al. Sampling Host-Seeking Anthropophilic Mosquito Vectors in West Africa: Comparisons of an Active Human-Baited Tent-Trap Against Gold Standard Methods. *The American Journal of Tropical Medicine and Hygiene*. 2015; 92: 415-21.
22. Cevallos V, Ponce P, Waggoner JJ, et al. Zika and Chikungunya virus detection in naturally infected *Aedes aegypti* in Ecuador. *Acta Tropica*. 2018; 177: 74-80.
23. Das AK, Harries AD, Hinderaker SG, et al. Active and passive case detection strategies for the control of leishmaniasis in Bangladesh. *Public Health Action*. 2014; 4: 15-21.
24. Insaf TZ and Talbot T. Identifying areas at risk of low birth weight using spatial epidemiology: a small area surveillance study. *Preventive medicine*. 2016; 88: 108-14.
25. Midega JT, Smith DL, Olotu A, et al. Wind direction and proximity to larval sites determines malaria risk in Kilifi District in Kenya. *Nature Communications*. 2012; 3: 674.
26. Vollack K, Sodoudi S, Névir P, Müller K and Richter D. Influence of meteorological parameters during the preceding fall and winter on the questing activity of nymphal *Ixodes ricinus* ticks. *International Journal of Biometeorology*. 2017; 61: 1787-95.
27. Gething PW, Noor AM, Gikandi PW, et al. Improving Imperfect Data from Health Management Information Systems in Africa Using Space–Time Geostatistics. *PLOS Medicine*. 2006; 3: e271.
28. Dickersin K and Chalmers I. Recognizing, investigating and dealing with incomplete and biased reporting of clinical research: from Francis Bacon to the WHO. *Journal of the Royal Society of Medicine*. 2011; 104: 532-8.
29. Smyth RMD, Kirkham JJ, Jacoby A, Altman DG, Gamble C and Williamson PR. Frequency and reasons for outcome reporting bias in clinical trials: interviews with trialists. *BMJ*. 2011; 342.
30. Nemet GF and Bailey AJ. Distance and health care utilization among the rural elderly. *Social Science & Medicine*. 2000; 50: 1197-208.
31. Kiwanuka SN, Ekirapa EK, Peterson S, et al. Access to and utilisation of health services for the poor in Uganda: a systematic review of available evidence. *Transactions of the Royal Society of Tropical Medicine and Hygiene*. 2008; 102: 1067-74.
32. Schoeps A, Gabrysch S, Niamba L, Sié A and Becher H. The Effect of Distance to Health-Care Facilities on Childhood Mortality in Rural Burkina Faso. *American Journal of Epidemiology*. 2011; 173: 492-8.
33. Kizito J, Kayendeke M, Nabirye C, Staedke SG and Chandler CIR. Improving access to health care for malaria in Africa: a review of literature on what attracts patients. *Malaria Journal*. 2012; 11: 55.
34. Larson PS, Mathanga DP, Campbell CH and Wilson ML. Distance to health services influences insecticide-treated net possession and use among six to 59 month-old children in Malawi. *Malaria Journal*. 2012; 11: 18.
35. Oduro AR, Maya ET, Akazili J, Baiden F, Koram K and Bojang K. Monitoring malaria using health facility based surveys: challenges and limitations. *BMC Public Health*. 2016; 16.
36. Weiss DJ, Nelson A, Gibson HS, et al. A global map of travel time to cities to assess inequalities in accessibility in 2015. *Nature*. 2018; 553: 333.
37. Gaudart J, Touré O, Dessay N, et al. Modelling malaria incidence with environmental dependency in a locality of Sudanese savannah area, Mali. *Malaria Journal*. 2009; 8: 61.
38. Wayant NM, Maldonado D, de Arias AR, Cousiño B and Goodin DG. Correlation between normalized difference vegetation index and malaria in a subtropical rain forest undergoing rapid anthropogenic alteration. *Geospatial health*. 2010; 4: 179-90.

39. Palaniyandi M. The role of Remote Sensing and GIS for spatial prediction of vector-borne diseases transmission: A systematic review. *Journal of vector borne diseases*. 2012; 49: 197.
40. Parham PE, Pople D, Christiansen-Jucht C, Lindsay S, Hinsley W and Michael E. Modeling the role of environmental variables on the population dynamics of the malaria vector *Anopheles gambiae sensu stricto*. *Malaria Journal*. 2012; 11: 271.
41. Hirve S, Singh SP, Kumar N, et al. Effectiveness and Feasibility of Active and Passive Case Detection in the Visceral Leishmaniasis Elimination Initiative in India, Bangladesh, and Nepal. *The American Journal of Tropical Medicine and Hygiene*. 2010; 83: 507-11.
42. Kuznetsov VN, Grjibovski AM, Mariandyshev AO, Johansson E and Bjune GA. A comparison between passive and active case finding in TB control in the Arkhangelsk region. *International Journal of Circumpolar Health*. 2014; 73: 23515.
43. Zhou G, Afrane YA, Malla S, Githeko AK and Yan G. Active case surveillance, passive case surveillance and asymptomatic malaria parasite screening illustrate different age distribution, spatial clustering and seasonality in western Kenya. *Malaria journal*. 2015; 14: 41.
44. Pava Z, Handayuni I, Trianty L, et al. Passively versus Actively Detected Malaria: Similar Genetic Diversity but Different Complexity of Infection. *The American Journal of Tropical Medicine and Hygiene*. 2017; 97: 1788-96.
45. Illian J, Penttinen A, Stoyan H and Stoyan D. *Statistical analysis and modelling of spatial point patterns*. John Wiley & Sons, 2008.
46. Wiegand T and Moloney KA. *Handbook of spatial point-pattern analysis in ecology*. CRC Press, 2013.
47. Buckland ST. *Introduction to Distance Sampling: Estimating Abundance of Biological Populations*. Oxford: Oxford University Press, 2001.
48. Chandler RB and Royle JA. Spatially explicit models for inference about density in unmarked or partially marked populations. *The Annals of Applied Statistics*. 2013; 7: 936-54.
49. Ramsey DS, Caley PA and Robley A. Estimating population density from presence-absence data using a spatially explicit model. *The Journal of Wildlife Management*. 2015; 79: 491-9.
50. Bhatt S, Weiss DJ, Cameron E, et al. The effect of malaria control on *Plasmodium falciparum* in Africa between 2000 and 2015. *Nature*. 2015; 526: 207.
51. Gething PW, Casey DC, Weiss DJ, et al. Mapping *Plasmodium falciparum* Mortality in Africa between 1990 and 2015. *New England Journal of Medicine*. 2016; 375: 2435-45.
52. Protopopoff N, Van Bortel W, Speybroeck N, et al. Ranking Malaria Risk Factors to Guide Malaria Control Efforts in African Highlands. *PLOS ONE*. 2009; 4: e8022.
53. Ranson H, N'Guessan R, Lines J, Moiroux N, Nkuni Z and Corbel V. Pyrethroid resistance in African anopheline mosquitoes: what are the implications for malaria control? *Trends in Parasitology*. 2011; 27: 91-8.
54. Toé KH, Jones CM, N'Fale S, Ismail HM, Dabiré RK and Ranson H. Increased Pyrethroid Resistance in Malaria Vectors and Decreased Bed Net Effectiveness, Burkina Faso. *Emerging Infectious Diseases*. 2014; 20: 1691-6.
55. Ranson H and Lissenden N. Insecticide Resistance in African *Anopheles* Mosquitoes: A Worsening Situation that Needs Urgent Action to Maintain Malaria Control. *Trends in Parasitology*. 2016; 32: 187-96.
56. Badolo A, Traore A, Jones CM, et al. Three years of insecticide resistance monitoring in *Anopheles gambiae* in Burkina Faso: resistance on the rise? *Malaria Journal*. 2012; 11: 232.
57. Govella NJ, Chaki PP and Killeen GF. Entomological surveillance of behavioural resilience and resistance in residual malaria vector populations. *Malaria Journal*. 2013; 12: 124.
58. Guelbeogo WM, Sagnon NF, Liu F, Besansky NJ and Costantini C. Behavioural divergence of sympatric *Anopheles funestus* populations in Burkina Faso. *Malaria Journal*. 2014; 13: 65-.
59. Mayer SV, Tesh RB and Vasilakis N. The emergence of arthropod-borne viral diseases: A global prospective on dengue, chikungunya and Zika fevers. *Acta Tropica*. 2017; 166: 155-63.

60. Dantas-Torres F, Chomel BB and Otranto D. Ticks and tick-borne diseases: a One Health perspective. *Trends in Parasitology*. 2012; 28: 437-46.
61. Lai Y-S, Biedermann P, Ekpo UF, et al. Spatial distribution of schistosomiasis and treatment needs in sub-Saharan Africa: a systematic review and geostatistical analysis. *The Lancet Infectious Diseases*. 2015; 15: 927-40.
62. Nanyingi MO, Munyua P, Kiama SG, et al. A systematic review of Rift Valley Fever epidemiology 1931–2014. *Infection Ecology & Epidemiology*. 2015; 5: 28024.
63. Franco JR, Simarro PP, Diarra A and Jannin JG. Epidemiology of human African trypanosomiasis. *Clinical Epidemiology*. 2014; 6: 257-75.
64. Davis JK, Vincent G, Hildreth MB, Kightlinger L, Carlson C and Wimberly MC. Integrating Environmental Monitoring and Mosquito Surveillance to Predict Vector-borne Disease: Prospective Forecasts of a West Nile Virus Outbreak. *PLoS Currents*. 2017; 9: ecurrents.outbreaks.90e80717c4e67e1a830f17feeaf85de.
65. Lawson AB. *Bayesian disease mapping: hierarchical modeling in spatial epidemiology*. CRC press, 2013.
66. Bousema T, Okell L, Felger I and Drakeley C. Asymptomatic malaria infections: detectability, transmissibility and public health relevance. *Nature Reviews Microbiology*. 2014; 12: 833.
67. Ripley BD. Spatial statistics. 1981. *Hayward Wiley, New York*. 1981.
68. Kelsall J and Wakefield J. Modeling Spatial Variation in Disease Risk: A Geostatistical Approach. *Journal of the American Statistical Association*. 2002; 97: 692-701.
69. Diboulo E, Sié A and Vounatsou P. Assessing the effects of malaria interventions on the geographical distribution of parasitaemia risk in Burkina Faso. *Malaria Journal*. 2016; 15: 228.
70. Royle JA. N-mixture models for estimating population size from spatially replicated counts. *Biometrics*. 2004; 60: 108-15.
71. Kéry M and Royle JA. *Applied Hierarchical Modeling in Ecology: Analysis of distribution, abundance and species richness in R and BUGS: Volume 1: Prelude and Static Models*. Academic Press, 2015.
72. Liu J and Chen X-P. Relationship of remote sensing normalized differential vegetation index to Anopheles density and malaria incidence rate. *Biomedical and environmental sciences: BES*. 2006; 19: 130-2.
73. Krefis AC, Schwarz NG, Krüger A, et al. Modeling the Relationship between Precipitation and Malaria Incidence in Children from a Holoendemic Area in Ghana. *The American Journal of Tropical Medicine and Hygiene*. 2011; 84: 285-91.
74. Diboulo E, Sié A, Diadier DA, Voules DAK, Yé Y and Vounatsou P. Bayesian variable selection in modelling geographical heterogeneity in malaria transmission from sparse data: an application to Nouna Health and Demographic Surveillance System (HDSS) data, Burkina Faso. *Parasites & vectors*. 2015; 8: 118.
75. Srimath-Tirumula-Peddinti RCPK, Neelapu NRR and Sidagam N. Association of Climatic Variability, Vector Population and Malarial Disease in District of Visakhapatnam, India: A Modeling and Prediction Analysis. *PLoS ONE*. 2015; 10: e0128377.
76. Alout H, Roche B, Dabiré RK and Cohuet A. Consequences of insecticide resistance on malaria transmission. *PLOS Pathogens*. 2017; 13: e1006499.
77. Coleman M, Hemingway J, Gleave KA, Wiebe A, Gething PW and Moyes CL. Developing global maps of insecticide resistance risk to improve vector control. *Malaria journal*. 2017; 16: 86.
78. Busetto L and Ranghetti L. MODISTsp: An R package for automatic preprocessing of MODIS Land Products time series. *Computers & Geosciences*. 2016; 97: 40-8.
79. Silverman BW. *Density estimation for statistics and data analysis*. CRC press, 1986.
80. Gitzen RA and Millspaugh JJ. Comparison of least-squares cross-validation bandwidth options for kernel home-range estimation. *Wildlife Society Bulletin*. 2003: 823-31.

81. Matowo N, Munhenga G, Tanner M, et al. *Fine-scale spatial and temporal heterogeneities in insecticide resistance profiles of the malaria vector, Anopheles arabiensis in rural south-eastern Tanzania [version 1; referees: 2 approved]*. 2017.
82. QGIS Development Team. QGIS Geographic Information System. Open Source Geospatial Foundation Project. 2018.
83. R Development Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing, 2018.
84. Plummer M. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. *Proceedings of the 3rd international workshop on distributed statistical computing*. Vienna, Austria, 2003, p. 125.
85. Plummer MS, Alexey Denwood, Matt. rjags: Bayesian Graphical Models using MCMC. Version 4.6. Downloaded from <https://cran.r-project.org/web/packages/rjags/index.html>. 2016.
86. World Health Organization. *Malaria entomology and vector control*. World Health Organization, 2013.
87. Corbel V and N'Guessan R. Distribution, mechanisms, impact and management of insecticide resistance in malaria vectors: a pragmatic review. *Anopheles mosquitoes-New insights into malaria vectors*. InTech, 2013.
88. Robinson OJ, Ruiz-Gutierrez V and Fink D. Correcting for bias in distribution modelling for rare species using citizen science data. *Diversity and Distributions*. 2017: n/a-n/a.
89. Mtema Z, Chagalucha J, Cleaveland S, et al. Mobile Phones As Surveillance Tools: Implementing and Evaluating a Large-Scale Intersectoral Surveillance System for Rabies in Tanzania. *PLOS Medicine*. 2016; 13: e1002002.
90. Lindblade KA, Steinhardt L, Samuels A, Kachur SP and Slutsker L. The silent threat: asymptomatic parasitemia and malaria transmission. *Expert Review of Anti-infective Therapy*. 2013; 11: 623-39.
91. Laake JL and Borchers DL. Methods for incomplete detection at distance zero. *Advanced Distance Sampling*. Oxford University Press, 2004, p. 108-89.
92. Alegana VA, Atkinson PM, Lourenço C, et al. Advances in mapping malaria for elimination: fine resolution modelling of Plasmodium falciparum incidence. *Scientific Reports*. 2016; 6: 29628.
93. Rue H, Martino S and Chopin N. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2009; 71: 319-92.
94. Blangiardo M and Cameletti M. *Spatial and spatio-temporal Bayesian models with R-INLA*. John Wiley & Sons, 2015.

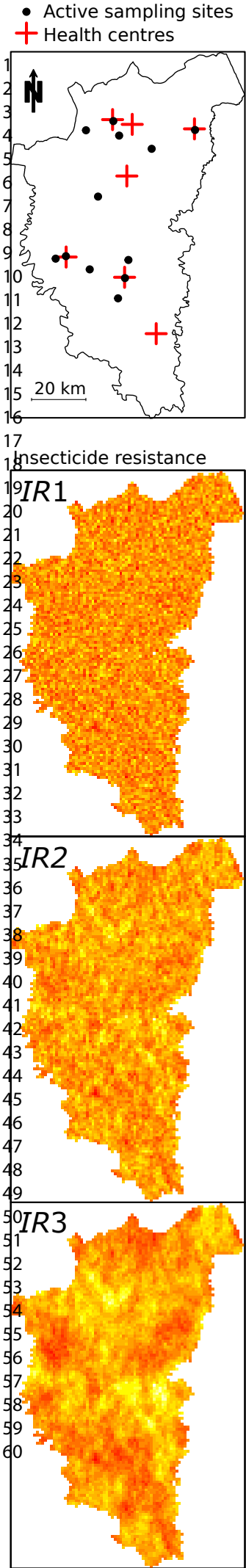
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

FIGURE CAPTIONS

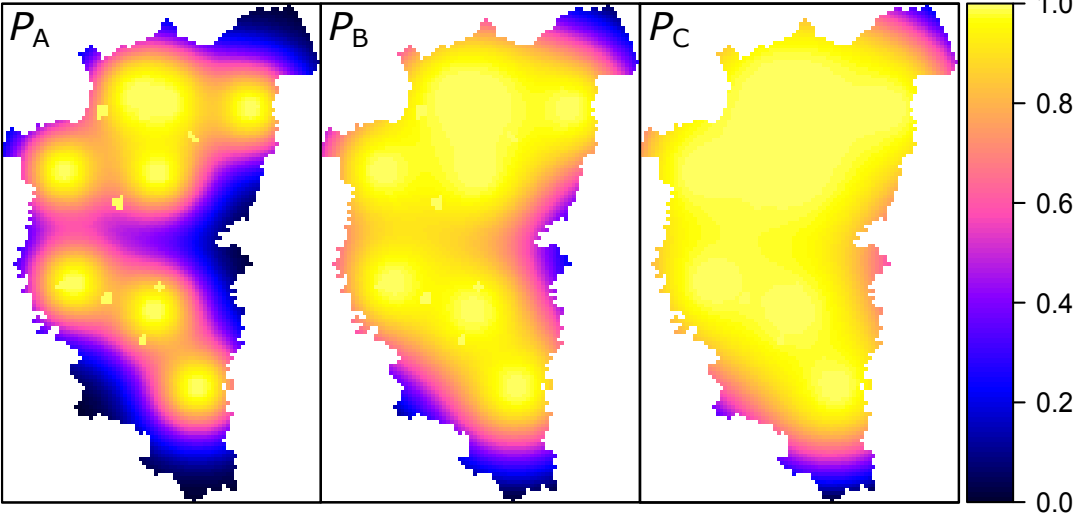
Figure 1 – Location of active sampling sites, simulation of reported malaria reported incidence (I_{1A} , ..., I_{3C}) under 3 scenarios of insecticide resistance ($IR1$, $IR2$, $IR3$) and 3 scenarios of reporting probability as a function of distance from health centres (P_A , P_B , P_C).

Figure 2 – Visual summary of results of the three Bayesian models of reported malaria incidence (I) under different simulated scenarios of insecticide resistance spatial patterns (IR) and probability of reporting at health centres (P). Model 1 used only active sampling data from some localised surveys, model 2 only passive case detections at health centres, model 3 combined both data sources together. The colour scale refers to the absolute values of the relative bias between the simulated coefficients of the variables involved in the biological process (1 to 6), or the shape parameter of the detection function (7), and the estimate of the same coefficient obtained by the mean of Markov Chain Monte Carlo (MCMC) posterior distributions. (•) indicates that the simulated coefficient is within the corresponding 95% posterior credible interval, (x) indicates that it falls outside.

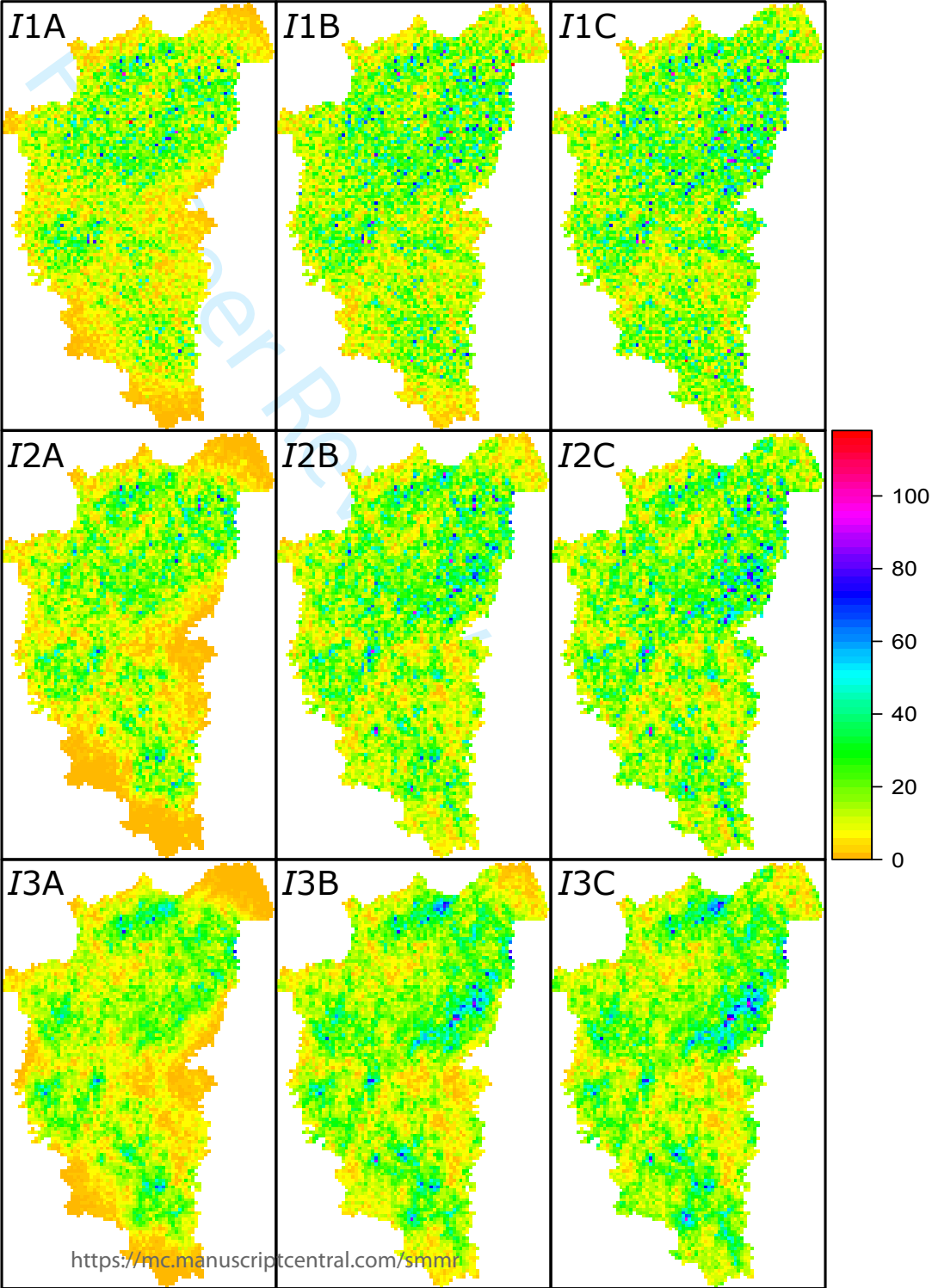
Figure 3 – Reconstructions of a simulated scenario of a) malaria incidence and b) insecticide resistance using Bayesian models. Figure refers to scenario 3B (see Fig. 1), with a high level of insecticide resistance spatial autocorrelation and an intermediate shape of the detection function. Model 1 used only active sampling data from a small set of localised surveys, model 2 only passive case detections at health centres, model 3 combined both data sources together.



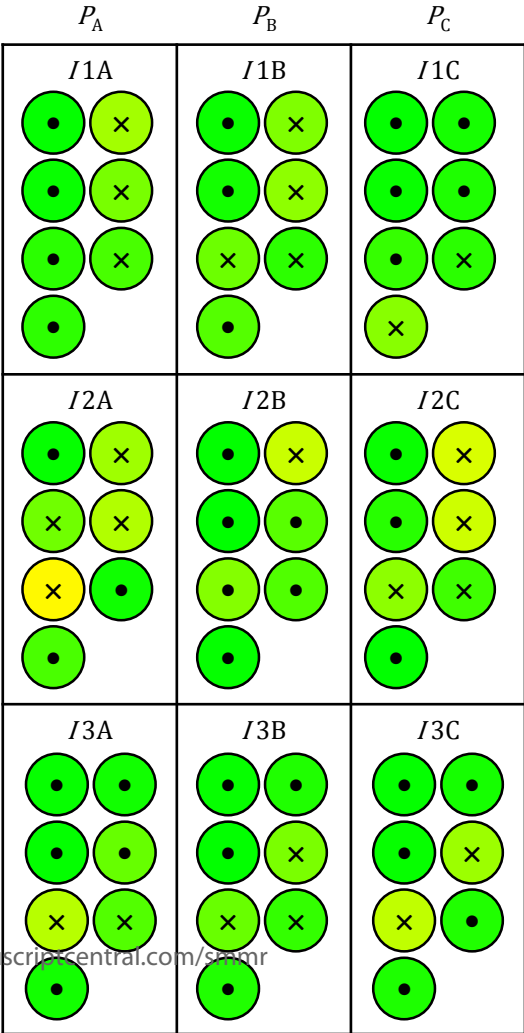
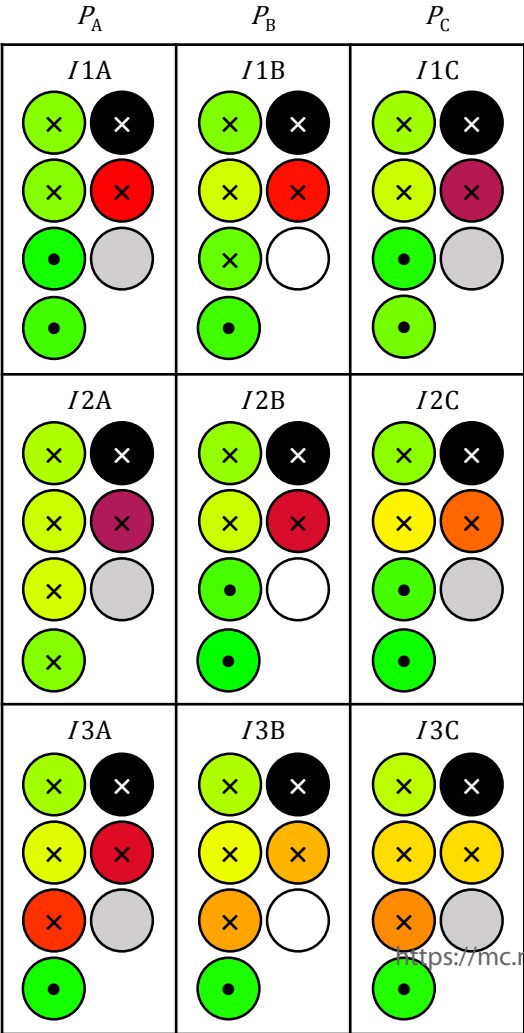
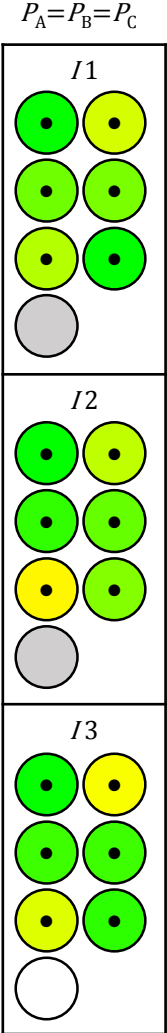
Probability of case reporting



Reported malaria incidence



1
2
3
4
5
6/IR 1
7
8
9
10
11
12
13
14
15
16/IR 2
17
18
19
20
21
22
23
24
25
26
27/IR 3
28
29
30
31
32



Predictors

1

2

3

4

5

6

7

1. Intercept

2. HUM

3. NDVI

4. RAIN

5. TEMP

6. IR

7. σ

x

Outside CI

•

Inside CI

|RB|

0.0

0.5

1.0

>1.0

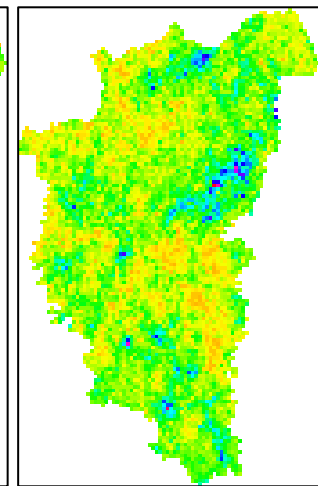
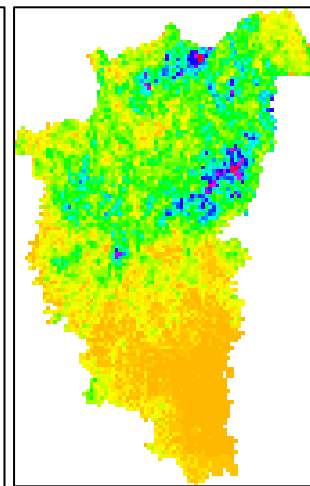
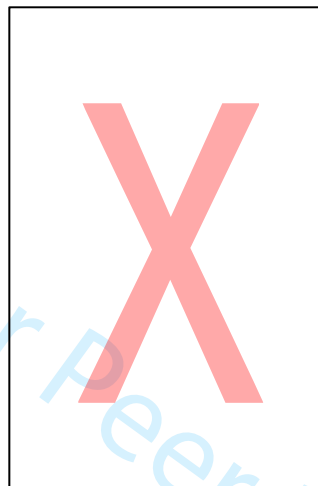
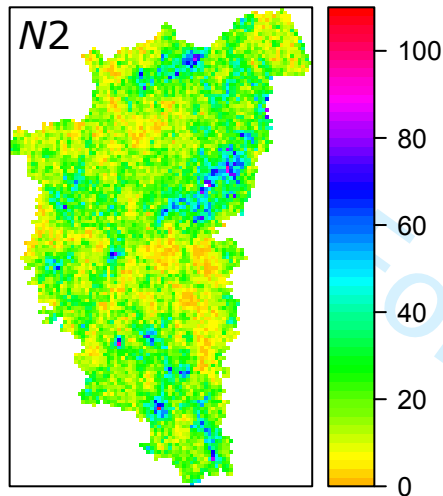
NA

a) Simulated

Model 1

Model 2

Model 3

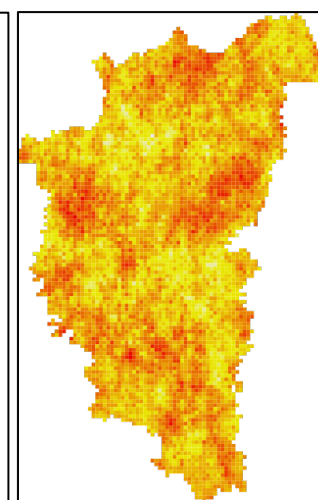
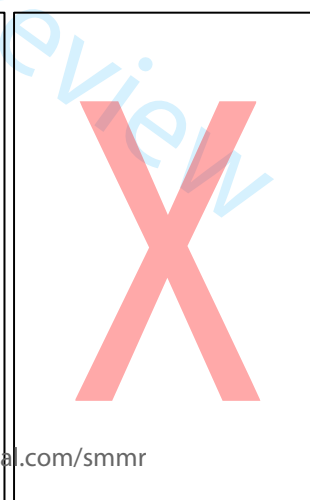
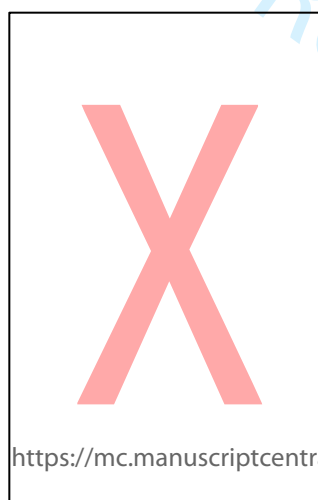
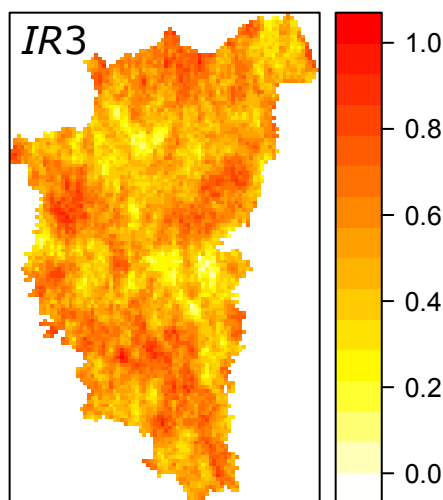


b) Simulated

Model 1

Model 2

Model 3



Supplementary material S1

Nelli L., Ferguson H.M., Matthiopoulos J.

Achieving depth and breadth in spatial models of vector-borne diseases: an integrated framework for active survey and passive surveillance data.

The following R scripts are used to generate the simulated dataset and to run the three models presented in the paper. Note that in the paper we presented 9 different scenarios, whereas here we are simulating only 1 scenario with an intermediate level of spatial autocorrelation of insecticide resistance ($ro=0.7$, see L 57) and an intermediate scenario of detectability ($\sigma=15$, see L 126). Also, note that the dataset provided as supplementary material is a subset of the entire data that we analysed in the paper, therefore some difference in the final results might be expected. For constant updates on this scripts, please visit <https://lucanelli.wordpress.com/r-codes>.

R code for data simulation and models

```
library(rgdal)
library(mvtnorm)
library(raster)
library(rgeos)
library(R2jags)
library(jagstools)

#read the grid in shapefile format. Please find it as paper's supplementary material
"grid.shp"
grid<-readOGR("your.folder.path","grid", GDAL1_integer64_policy=T)

# set the distance from clinics=0 in active sampling points (to set probability of reporting
as 1)
grid$dist_1[grid$sampling==1] <- 0
grid$dist_2[grid$sampling==1] <- 0
grid$dist_3[grid$sampling==1] <- 0
grid$dist_4[grid$sampling==1] <- 0
grid$dist_5[grid$sampling==1] <- 0
grid$dist_6[grid$sampling==1] <- 0
grid$dist_7[grid$sampling==1] <- 0
grid$dist_8[grid$sampling==1] <- 0

# #####
# ---- - simulate Insecticide Resistance - ----
# #####

# function to make a distance matrix (side*side) 2D array
dist.matrix <- function(side)
{
  row.coords <- rep(1:side, times=side)
  col.coords <- rep(1:side, each=side)
  row.col <- data.frame(row.coords, col.coords)
  D <- dist(row.col, method="euclidean", diag=TRUE, upper=TRUE)
  D <- as.matrix(D)
  return(D)
}

# function to simulate an autocorrelated surface, with exponential decay given by ro
cor.surface <- function(side, global.mu, ro)
{
  D <- dist.matrix(side)
  # scaling the distance matrix by the exponential decay
  SIGMA <- exp(-ro*D)
  mu <- rep(global.mu, times=side*side)
  # sampling from the multivariate normal distribution
  M <- matrix(nrow=side, ncol=side)
```

```

1
2
3 50 M[] <- rmvnorm(1, mu, SIGMA)
4 51 return(M)
5 52 }
6 53
7 54 # parameters
8 55 side <- max(c(grid@bbox[1,2]-grid@bbox[1,1],grid@bbox[2,2]-grid@bbox[2,1]))/1000 # Arena
9 56 dimension (take the maximum extension of the bounding box of the grid, in Km)
10 57 global.mu <- 0
11 58 ro <- 0.7 #note that here I'm simulating only the intermediate scenario.
12 59
13 60 # simulating the autocorrelated raster
14 61 set.seed(1)
15 62 ir <- cor.surface(side = side, ro = ro, global.mu = global.mu) #this may take a lot (up to few
16 63 hours, depending on the computation power)
17 64 image(ir) # have a look
18 65
19 66 # now transform it into a raster and assign it to the cells
20 67 ir.raster <- raster(ir)
21 68 extent(ir.raster) <- (grid@bbox)
22 69 plot(ir.raster)
23 70 summary(ir.raster)
24 71 grid$IR<-(extract(ir.raster, gCentroid(grid,byid=TRUE), na.rm = T, small = T, df = T))$layer
25 72 summary(grid$IR)
26 73
27 74 #####
28 75 #### - Simulate the dataset of malaria cases - ####
29 76 #####
30 77
31 78 #scaling environmental variables
32 79 scale2 <- function(x) {
33 80   sdx <- sqrt(var(x))
34 81   meanx <- mean(x)
35 82   return((x - meanx)/sdx)
36 83 }
37 84
38 85 grid$NDVI<-scale2(grid$NDVI)
39 86 grid$RAIN<-scale2(as.integer(grid$RAIN))
40 87 grid$TEMP<-scale2(as.integer(grid$TEMP))
41 88 grid$HUM<-scale2(grid$HUM)
42 89
43 90
44 91 # set the simulated coefficients and create the matrix of coefficients 'a'
45 92 sim.intercept <- 2.90
46 93 sim.beta.NDVI <- 0.30
47 94 sim.beta.RAIN <- 0.20
48 95 sim.beta.TEMP <- 0.25
49 96 sim.beta.HUM <- 0.20
50 97 sim.beta.IR <- 0.50
51 98
52 99
53 100 a<-rbind(sim.intercept,
54 101          sim.beta.NDVI,
55 102          sim.beta.RAIN,
56 103          sim.beta.TEMP,
57 104          sim.beta.HUM,
58 105          sim.beta.IR)
59 106
60 107 #Create a the covariate matrix 'x'
61 108 x<- matrix(c(rep(1,length(grid)),
62 109            grid$NDVI,
63 110            grid$RAIN,
64 111            grid$TEMP,
65 112            grid$HUM,
66 113            grid$IR),
67 114            nrow=length(grid))
68 115
69 116 #generate a poisson rates L
70 117 grid$L<- as.vector(x %*% a)
71 118
72 119 #generate true incidence N (number of malaria cases under the biological process)
73 120 grid$lambda<-exp(grid$L)
74 121 set.seed(1)
75 122 grid$N<-rpois(length(grid),grid$lambda)
76 123
77 124 # set probability of reporting (Pd) at each clinic, according to a half normal function with
78 125 sigma=15
79 126 HN<-function (x,si) {exp(-(x^2)/(2*sigma^2))}

```

```

1
2
3 127 sigma<-15
4 128
5 129 grid$Pd1<-HN(grid$dist_1,si)
6 130 grid$Pd2<-HN(grid$dist_2,si)
7 131 grid$Pd3<-HN(grid$dist_3,si)
8 132 grid$Pd4<-HN(grid$dist_4,si)
9 133 grid$Pd5<-HN(grid$dist_5,si)
10 134 grid$Pd6<-HN(grid$dist_6,si)
11 135 grid$Pd7<-HN(grid$dist_7,si)
12 136 grid$Pd8<-HN(grid$dist_8,si)
13 137
14 138 # set the overall probability of not reporting any case at all (Q)
15 139 grid$Q<-(1-grid$Pd1)*
16 140 (1-grid$Pd2)*
17 141 (1-grid$Pd3)*
18 142 (1-grid$Pd4)*
19 143 (1-grid$Pd5)*
20 144 (1-grid$Pd6)*
21 145 (1-grid$Pd7)*
22 146 (1-grid$Pd8)
23 147
24 148 # generate the reported incidence (I) under the observation process
25 149 n.clinics<-8
26 150 grid$I<-matrix(nrow = length(grid), ncol = n.clinics+1)
27 151
28 152 for (i in 1:length(grid)) {
29 153   set.seed(1)
30 154   grid$I[i,]<-rmultinom(1, grid$N[i], cbind(grid$Pd1[i],
31 155                                             grid$Pd2[i],
32 156                                             grid$Pd3[i],
33 157                                             grid$Pd4[i],
34 158                                             grid$Pd5[i],
35 159                                             grid$Pd6[i],
36 160                                             grid$Pd7[i],
37 161                                             grid$Pd8[i],
38 162                                             grid$Q[i]))
39 163 }
40 164
41 165 # total reported cases at each cell
42 166 grid$total.reported.cases<-rowSums(grid$I[,1:n.clinics])
43 167
44 168 # reported cases from clinics only (set NA in active sampling cells)
45 169 grid$incidence.clinics<-grid$I
46 170 grid$incidence.clinics[grid$sampling==1]<-NA
47 171
48 172
49 173 #####
50 174 ##### Model 1 - active sampling only #####
51 175 #####
52 176
53 177 #subset the data, only active sampling sites
54 178 grid.sub<-subset(grid, sampling==1)
55 179
56 180 #Jags model
57 181 model.1<-function() {
58 182   # Priors
59 183   alpha ~ dnorm(0,0.001)
60 184   beta.NDVI ~ dnorm(0,0.001)
61 185   beta.RAIN ~ dnorm(0,0.001)
62 186   beta.TEMP ~ dnorm(0,0.001)
63 187   beta.HUM ~ dnorm(0,0.001)
64 188   beta.IR ~ dnorm(0,0.001)
65 189   si ~ dgamma(0.01, 0.01)
66 190
67 191   for (i in 1:n) {
68 192
69 193     count[i]~dbin(1,raw.count[i])
70 194
71 195     raw.count[i]~dpois(lambda[i])
72 196
73 197     log(lambda[i]) <- alpha +
74 198       beta.NDVI*NDVI[i] +
75 199       beta.RAIN*RAIN[i] +
76 200       beta.TEMP*TEMP[i]+
77 201       beta.HUM*HUM[i]+
78 202       beta.IR*IR[i]+eps[i]
79 203

```

```

204     eps[i]~dnorm(0, si)
205   }
206 }
207 }
208
209
210 params<-c("alpha","beta.NDVI","beta.RAIN","beta.TEMP","beta.HUM","beta.IR")
211 n.iterations<-2000
212
213 jags.data.survey <- list (count=grid.sub$total.reported.cases,
214   NDVI=grid.sub$NDVI,
215   RAIN=grid.sub$RAIN,
216   TEMP=grid.sub$TEMP,
217   HUM=grid.sub$HUM,
218   IR=grid.sub$IR,
219   n=length(grid.sub))
220
221
222 jags.out.survey<-jags(data=jags.data.survey,
223   model.file=model.1,
224   n.chains= 3,
225   n.iter=n.iterations,
226   parameters.to.save=params)
227
228
229 #####
230 ##### Model 2 - passive case detection #####
231 #####
232
233
234 model.2<-function() {
235   # Priors
236   alpha      ~ dnorm(0,0.0001)
237   beta.NDVI   ~ dnorm(0,0.0001)
238   beta.RAIN   ~ dnorm(0,0.0001)
239   beta.TEMP   ~ dnorm(0,0.0001)
240   beta.HUM    ~ dnorm(0,0.0001)
241   sigma      ~ dunif(0,100)
242   si          ~ dgamma(0.01, 0.01)
243
244   # Likelihood
245
246   for (i in 1:n) {
247
248     counts[i,1:8]~dmulti(probs[i,1:8],reported.cases[i])
249
250     probs[i,1]<-exp(-(dist_1[i]*dist_1[i])/(2*sigma*sigma))
251     probs[i,2]<-exp(-(dist_2[i]*dist_2[i])/(2*sigma*sigma))
252     probs[i,3]<-exp(-(dist_3[i]*dist_3[i])/(2*sigma*sigma))
253     probs[i,4]<-exp(-(dist_4[i]*dist_4[i])/(2*sigma*sigma))
254     probs[i,5]<-exp(-(dist_5[i]*dist_5[i])/(2*sigma*sigma))
255     probs[i,6]<-exp(-(dist_6[i]*dist_6[i])/(2*sigma*sigma))
256     probs[i,7]<-exp(-(dist_7[i]*dist_7[i])/(2*sigma*sigma))
257     probs[i,8]<-exp(-(dist_8[i]*dist_8[i])/(2*sigma*sigma))
258
259     reported.cases[i]~dbinom(overall_prob[i], true.incidence[i])
260
261     overall_prob[i]<-1-Q[i]
262
263     Q[i] <- q1[i]*q2[i]*q3[i]*q4[i]*q5[i]*q6[i]*q7[i]*q8[i]
264     q1[i] <- 1-probs[i,1]
265     q2[i] <- 1-probs[i,2]
266     q3[i] <- 1-probs[i,3]
267     q4[i] <- 1-probs[i,4]
268     q5[i] <- 1-probs[i,5]
269     q6[i] <- 1-probs[i,6]
270     q7[i] <- 1-probs[i,7]
271     q8[i] <- 1-probs[i,8]
272
273     true.incidence[i]~dpois(lambda[i])
274
275     log(lambda[i]) <- alpha + beta.NDVI*NDVI[i] +
276       beta.RAIN*RAIN[i] +
277       beta.TEMP*TEMP[i] +
278       beta.HUM*HUM[i] + eps[i]
279
280

```

```

1
2
3 281     eps[i]~dnorm(0, si)
4 282   }
5 283 }
6 284
7 285
8 286 params<-c("alpha","beta.NDVI","beta.RAIN","beta.TEMP","beta.HUM", "sigma","true.incidence")
9 287
10 288 inits <- function(){list (true.incidence = rowSums(counts[,1:8]))}
11 289
12 290 n.iterations<-2000
13 291
14 292
15 293 jags.data.clinics <-list (counts=grid$incidence.clinics[,1:8],
16 294                          NDVI=grid$NDVI,
17 295                          RAIN=grid$RAIN,
18 296                          TEMP=grid$TEMP,
19 297                          HUM=grid$HUM,
20 298                          dist_1=grid$dist_1,
21 299                          dist_2=grid$dist_2,
22 300                          dist_3=grid$dist_3,
23 301                          dist_4=grid$dist_4,
24 302                          dist_5=grid$dist_5,
25 303                          dist_6=grid$dist_6,
26 304                          dist_7=grid$dist_7,
27 305                          dist_8=grid$dist_8,
28 306                          n=length(grid))
29 307
30 308 jags.out.clinics<-jags(data=jags.data.clinics,
31 309                       model.file=model.2,
32 310                       inits=inits,
33 311                       n.chains=3,
34 312                       n.iter=n.iterations,
35 313                       parameters.to.save=params)
36 314
37 315 #####
38 316 ##### Model 3 - active and passive case together #####
39 317 #####
40 318
41 319
42 320 DM<-gDistance(gCentroid(grid,byid=TRUE),byid=T)/1000 # calculate distances for covariance
43 321 matrix
44 322
45 323 model.3<-function() {
46 324   # Priors
47 325   alpha      ~ dnorm(0,0.0001)
48 326   beta.NDVI   ~ dnorm(0,0.0001)
49 327   beta.RAIN   ~ dnorm(0,0.0001)
50 328   beta.TEMP   ~ dnorm(0,0.0001)
51 329   beta.HUM    ~ dnorm(0,0.0001)
52 330   beta.IR     ~ dnorm(0,0.0001)
53 331   si          ~ dgamma(0.01, 0.01)
54 332   sigma       ~ dunif(0, 100)
55 333   ro          ~ dgamma(0.1,0.1)
56 334   global.mu   ~ dnorm(0,0.0001)
57 335
58 336
59 337   # Likelihood
60 338
61 339   for (i in 1:n) {
62 340
63 341     counts[i,1:8]~dmulti(probs[i,1:8],reported.cases[i])
64 342
65 343     probs[i,1]<-exp(-(dist_1[i]*dist_1[i])/(2*sigma*sigma))
66 344     probs[i,2]<-exp(-(dist_2[i]*dist_2[i])/(2*sigma*sigma))
67 345     probs[i,3]<-exp(-(dist_3[i]*dist_3[i])/(2*sigma*sigma))
68 346     probs[i,4]<-exp(-(dist_4[i]*dist_4[i])/(2*sigma*sigma))
69 347     probs[i,5]<-exp(-(dist_5[i]*dist_5[i])/(2*sigma*sigma))
70 348     probs[i,6]<-exp(-(dist_6[i]*dist_6[i])/(2*sigma*sigma))
71 349     probs[i,7]<-exp(-(dist_7[i]*dist_7[i])/(2*sigma*sigma))
72 350     probs[i,8]<-exp(-(dist_8[i]*dist_8[i])/(2*sigma*sigma))
73 351
74 352     reported.cases[i]~dbinom(overall_prob[i], true.incidence[i])
75 353
76 354     overall_prob[i]<-1-Q[i]
77 355
78 356     Q[i] <- q1[i]*q2[i]*q3[i]*q4[i]*q5[i]*q6[i]*q7[i]*q8[i]
79 357     q1[i] <- 1-probs[i,1]

```

```

358 q2[i] <- 1-probs[i,2]
359 q3[i] <- 1-probs[i,3]
360 q4[i] <- 1-probs[i,4]
361 q5[i] <- 1-probs[i,5]
362 q6[i] <- 1-probs[i,6]
363 q7[i] <- 1-probs[i,7]
364 q8[i] <- 1-probs[i,8]
365
366 true.incidence[i]~dpois(lambda[i])
367
368 log(lambda[i]) <- alpha + beta.NDVI*NDVI[i] +
369   beta.IR*IR[i]+
370   beta.RAIN*RAIN[i] +
371   beta.TEMP*TEMP[i] +
372   beta.HUM*HUM[i] + eps[i]
373
374 eps[i]~dnorm(0, si)
375
376
377 for(j in 1:n)
378 {
379   # turning the distance matrix to covariance matrix
380   C.w[i,j] <- exp(-ro*D[i,j])
381 }
382
383 # turning covariances into precisions
384 P.w[i,i] <- inverse(C.w[i,i])
385 mu[i] <- global.mu
386 IR[i] ~ dmnorm(mu[i], P.w[i,i])
387
388 }
389 }
390
391
392
393 n.iterations<-20000
394
395 params<-c("alpha","beta.NDVI","beta.RAIN","beta.HUM", "beta.TEMP", "beta.IR", "sigma","ro",
396 "global.mu","true.incidence","IR")
397
398 grid$IR.NA<-grid$IR
399 grid$IR.NA[grid$sampling==0]<-NA #set NA in all the cells but those of sampling points
400 counts=grid$I[,1:8]
401 inits <- function(){list (true.incidence = rowSums(counts[,1:8]))}
402
403 index<-c(1:nrow(grid))[-which(grid$sampling==1)] # create an index for NAs cells (that is,
404 all but the 12 sampling stations)
405
406 jags.data.both <- list (counts=grid$I[,1:8],
407   NDVI=grid$NDVI,
408   RAIN=grid$RAIN,
409   TEMP=grid$TEMP,
410   HUM=grid$HUM,
411   dist_1=grid$dist_1,
412   dist_2=grid$dist_2,
413   dist_3=grid$dist_3,
414   dist_4=grid$dist_4,
415   dist_5=grid$dist_5,
416   dist_6=grid$dist_6,
417   dist_7=grid$dist_7,
418   dist_8=grid$dist_8,
419   IR=grid$IR.NA,
420   n=length(grid),
421   nn=length(index),
422   index=index,
423   D=DM)
424
425 jags.out.both<-jags(data=jags.data.both, #this will take a lot.
426   model.file=model.3,
427   n.chains=3,
428   inits=inits,
429   n.iter=n.iterations,
430   # n.thin=1,
431   parameters.to.save=params)
432
433
434

```

```

1
2
3 435
4 436 #####
5 437 ## RESULTS - MODEL 1 ##
6 438 #####
7 439
8 440 # results of model 1
9 441
10 442 mu.model1<-jagsresults(x=jags.out.survey, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
11 443 "beta.TEMP","beta.IR", "sigma"))[,1]
12 444 sd.model1<-jagsresults(x=jags.out.survey, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
13 445 "beta.TEMP","beta.IR", "sigma"))[,2]
14 446 LCI.model1<-jagsresults(x=jags.out.survey,
15 447 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP","beta.IR", "sigma"))[,3]
16 448 UCI.model1<-jagsresults(x=jags.out.survey,
17 449 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP","beta.IR", "sigma"))[,7]
18 450 Rhat.model1<-jagsresults(x=jags.out.survey,
19 451 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP","beta.IR", "sigma"))[,8]
20 452 n.eff.model1<-jagsresults(x=jags.out.survey,
21 453 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP","beta.IR", "sigma"))[,9]
22 454
23 455
24 456 res.model1<-data.frame(mu.vect=mu.model1,
25 457 sd.vect=sd.model1,
26 458 LCI=LCI.model1,
27 459 UCI=UCI.model1,
28 460 Rhat=Rhat.model1,
29 461 n.eff=n.eff.model1)
30 462
31 463 res.model1
32 464
33 465 #####
34 466 ## RESULTS - MODEL 2 ##
35 467 #####
36 468
37 469 # results of model 2
38 470
39 471
40 472 mu.model2<-jagsresults(x=jags.out.clinics,
41 473 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,1]
42 474 sd.model2<-jagsresults(x=jags.out.clinics,
43 475 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,2]
44 476 LCI.model2<-jagsresults(x=jags.out.clinics,
45 477 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,3]
46 478 UCI.model2<-jagsresults(x=jags.out.clinics,
47 479 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,7]
48 480 Rhat.model2<-jagsresults(x=jags.out.clinics,
49 481 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,8]
50 482 n.eff.model2<-jagsresults(x=jags.out.clinics,
51 483 params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP", "sigma"))[,9]
52 484
53 485
54 486 res.model2<-data.frame(mu.vect=mu.model2,
55 487 sd.vect=sd.model2,
56 488 LCI=LCI.model2,
57 489 UCI=UCI.model2,
58 490 Rhat=Rhat.model2,
59 491 n.eff=n.eff.model2)
60 492
61 493 res.model2
62 494
63 495 # reconstruction of true incidence from model 2
64 496
65 497 grid$rec.incid.model2<-jags.out.clinics$BUGSoutput$mean$true.incidence
66 498
67 499 # transform into a raster (for nicer plotting)
68 500 r<-raster()
69 501 extent(r) <- extent(grid)
70 502 rec.incid.model2.rast<-rasterize(grid,r,'rec.incid.model2')
71 503 N.rast<-rasterize(grid,r,'N')
72 504
73 505 par(mfrow=c(1, 3))
74 506 plot(N.rast, main ="True Incidence")
75 507 plot(rec.incid.model2.rast, main ="Reconstructed Incidence")
76 508 plot(grid$rec.incid.model2~grid$N, main="True vs Reconstructed Incidence")
77 509 par(mfrow=c(1, 1))
78 510
79 511

```



```

#####
## RESULTS - MODEL 3 ##
#####

# results of model 3

mu.model3<-jagsresults(x=jags.out.both, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
"beta.TEMP","beta.IR", "sigma"))[,1]
sd.model3<-jagsresults(x=jags.out.both, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
"beta.TEMP","beta.IR", "sigma"))[,2]
LCI.model3<-jagsresults(x=jags.out.both, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
"beta.TEMP","beta.IR", "sigma"))[,3]
UCI.model3<-jagsresults(x=jags.out.both, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
"beta.TEMP","beta.IR", "sigma"))[,7]
Rhat.model3<-jagsresults(x=jags.out.both, params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN",
"beta.TEMP","beta.IR", "sigma"))[,8]
n.eff.model3<-jagsresults(x=jags.out.both,
params=c("alpha","beta.HUM","beta.NDVI","beta.RAIN", "beta.TEMP","beta.IR", "sigma"))[,9]

res.model3<-data.frame(mu.vect=mu.model3,
sd.vect=sd.model3,
LCI=LCI.model3,
UCI=UCI.model3,
Rhat=Rhat.model3,
n.eff=n.eff.model3)

res.model3

# reconstruction of true incidence from model 3

grid$rec.incid.model3<-jags.out.both$BUGSoutput$mean$true.incidence

# transform into a raster (for nicer plotting)
r<-raster()
extent(r) <- extent(grid)
rec.incid.model3.rast<-rasterize(grid,r,'rec.incid.model3')

par(mfrow=c(1, 3))
plot(N.rast, main ="True Incidence")
plot(rec.incid.model3.rast, main ="Reconstructed Incidence")
plot(grid$rec.incid.model3~grid$N, main="True vs Reconstructed Incidence")
par(mfrow=c(1, 1))

# reconstruction of insecticide resistance from model 3

grid$rec.InsRes.model3<-jags.out.both$BUGSoutput$mean$IR

# transform into a raster (for nicer plotting)
r<-raster()
extent(r) <- extent(grid)
rec.InsRes.model3.rast<-rasterize(grid,r,'rec.InsRes.model3')

par(mfrow=c(1, 3))
plot(ir.raster, col=rev(heat.colors(255)), main ="True insecticide resistance")
plot(rec.InsRes.model3.rast,col=rev(heat.colors(255)), main ="Reconstructed insecticide
resistance")
plot(grid$rec.InsRes.model3~grid$IR, main="True vs Reconstructed insecticide resistance")
par(mfrow=c(1, 1))

```

Supplementary material S2

Nelli L., Ferguson H.M., Matthiopoulos J.

Achieving depth and breadth in spatial models of vector-borne diseases: an integrated framework for active survey and passive surveillance data.

S2.1 – Result of Bayesian models of reported malaria incidence under different scenarios of insecticide resistance patterns. The table shows results of model 1, which considered only active sampling data from some localised surveys. θ : simulated coefficient, $\hat{\theta}$: mean of posterior distribution, CI: credible interval, RB: relative bias.

Scenario $I_{1A}=I_{1B}=I_{1C}$			
Variable	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.89 (2.82, 2.95)	-0.01
HUM	0.20	0.24 (0.14, 0.34)	0.21
NDVI	0.30	0.27 (0.20, 0.33)	-0.11
RAIN	0.20	0.22 (0.14, 0.31)	0.12
TEMP	0.25	0.21 (0.08, 0.34)	-0.18
IR	0.50	0.49 (0.44, 0.55)	-0.01
Scenario $I_{2A}=I_{2B}=I_{2C}$			
Variable	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.90 (2.83, 2.96)	0.00
HUM	0.20	0.24 (0.14, 0.34)	0.19
NDVI	0.30	0.29 (0.21, 0.36)	-0.05
RAIN	0.20	0.18 (0.09, 0.27)	-0.10
TEMP	0.25	0.19 (0.05, 0.32)	-0.26
IR	0.50	0.44 (0.37, 0.50)	-0.13
Scenario $I_{3A}=I_{3B}=I_{3C}$			
Variable	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.86 (2.79, 2.92)	-0.01
HUM	0.20	0.15 (0.04, 0.26)	-0.26
NDVI	0.30	0.32 (0.24, 0.39)	0.06
RAIN	0.20	0.21 (0.11, 0.32)	0.06
TEMP	0.25	0.30 (0.15, 0.47)	0.22
IR	0.50	0.52 (0.44, 0.60)	0.04

S2.2 – Result of Bayesian models of reported malaria incidence under different scenarios of insecticide resistance patterns and detectability at health centres. The table shows results of model 2, which considered only passive case detections at health centres. σ : shape parameter of half-normal detection function, θ : simulated coefficient, $\hat{\theta}$: mean of posterior distribution, CI: credible interval, RB: relative bias.

Scenario I_{1A}				Scenario I_{1B}			Scenario I_{1C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.53 (2.51, 2.55)	-0.13	2.90	2.51 (2.49, 2.53)	-0.13	2.90	2.48 (2.46, 2.50)	-0.15
HUM	0.20	1.01 (0.99, 1.04)	4.07	0.20	1.07 (1.04, 1.09)	4.33	0.20	1.08 (1.05, 1.10)	4.38
NDVI	0.30	0.26 (0.24, 0.28)	-0.13	0.30	0.24 (0.22, 0.26)	-0.20	0.30	0.24 (0.22, 0.26)	-0.19
RAIN	0.20	0.30 (0.28, 0.33)	0.51	0.20	0.30 (0.27, 0.32)	0.48	0.20	0.33 (0.30, 0.35)	0.63
TEMP	0.25	0.25 (0.22, 0.27)	-0.02	0.25	0.23 (0.20, 0.25)	-0.10	0.25	0.24 (0.22, 0.26)	-0.03
σ	10.00	10.35 (9.99, 10.39)	0.03	15.00	14.14 (14.00, 14.19)	-0.06	20.00	21.74 (19.64, 21.83)	0.09

Scenario I_{2A}				Scenario I_{2B}			Scenario I_{2C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.50 (2.48, 2.52)	-0.14	2.90	2.50 (2.48, 2.52)	-0.14	2.90	2.48 (2.46, 2.50)	-0.14
HUM	0.20	1.10 (1.07, 1.13)	4.51	0.20	1.04 (1.02, 1.07)	4.21	0.20	1.04 (1.02, 1.07)	4.20
NDVI	0.30	0.24 (0.22, 0.26)	-0.19	0.30	0.24 (0.22, 0.26)	-0.19	0.30	0.22 (0.02, 0.24)	-0.26
RAIN	0.20	0.33 (0.30, 0.35)	0.63	0.20	0.31 (0.29, 0.34)	0.57	0.20	0.28 (0.26, 0.30)	0.40
TEMP	0.25	0.20 (0.17, 0.22)	-0.21	0.25	0.23 (0.21, 0.26)	-0.07	0.25	0.23 (0.21, 0.26)	-0.07
σ	10.00	9.46 (9.42, 9.50)	-0.05	15.00	14.93 (14.87, 15.02)	0.00	20.00	20.36 (19.27, 20.45)	0.02

Scenario I_{3A}				Scenario I_{3B}			Scenario I_{3C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.42 (2.40, 2.45)	-0.16	2.90	2.42 (2.40, 2.43)	-0.17	2.90	2.39 (2.37, 2.41)	-0.18
HUM	0.20	1.05 (1.02, 1.08)	4.25	0.20	1.05 (1.02, 1.07)	4.25	0.20	1.04 (1.02, 1.07)	4.22
NDVI	0.30	0.24 (0.22, 0.26)	-0.19	0.30	0.23 (0.21, 0.25)	-0.23	0.30	0.22 (0.20, 0.24)	-0.28
RAIN	0.20	0.31 (0.29, 0.34)	0.56	0.20	0.26 (0.24, 0.29)	0.31	0.20	0.26 (0.23, 0.28)	0.29
TEMP	0.25	0.14 (0.11, 0.16)	-0.45	0.25	0.16 (0.14, 0.19)	-0.34	0.25	0.16 (0.14, 0.18)	-0.36
σ	10.00	9.84 (9.80, 10.01)	-0.02	15.00	14.59 (14.53, 15.65)	-0.03	20.00	19.32 (19.23, 20.41)	-0.03

S2.3 – Result of Bayesian models of reported malaria incidence under different scenarios of insecticide resistance patterns and detectability at health centres. The table shows results of model 3, which considered both active surveys and passive case detections at health centres. σ : shape parameter oh half-normal detection function, θ : simulated coefficient, $\hat{\theta}$: mean of posterior distribution, CI: credible interval, RB: relative bias.

Scenario I _{1A}				Scenario I _{1B}			Scenario I _{1C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.91 (2.87, 2.96)	0.00	2.90	2.92 (2.87, 2.96)	0.01	2.90	2.90 (2.85, 2.94)	0.00
HUM	0.20	0.17 (0.15, 0.19)	-0.16	0.20	0.17 (0.16, 0.19)	-0.13	0.20	0.20 (0.18, 0.21)	-0.02
NDVI	0.30	0.30 (0.29, 0.32)	0.01	0.30	0.31 (0.29, 0.32)	0.02	0.30	0.30 (0.29, 0.32)	0.01
RAIN	0.20	0.18 (0.16, 0.20)	-0.12	0.20	0.17 (0.15, 0.19)	-0.14	0.20	0.21 (0.19, 0.23)	0.03
TEMP	0.25	0.24 (0.22, 0.26)	-0.04	0.25	0.22 (0.20, 0.24)	-0.11	0.25	0.24 (0.22, 0.26)	-0.05
IR	0.50	0.47 (0.45, 0.48)	-0.07	0.50	0.48 (0.46, 0.49)	-0.04	0.50	0.48 (0.46, 0.49)	-0.04
σ	10.00	10.37 (10.00, 10.41)	0.04	15.00	14.80 (14.74, 15.85)	-0.05	20.00	22.58 (22.48, 22.67)	0.13

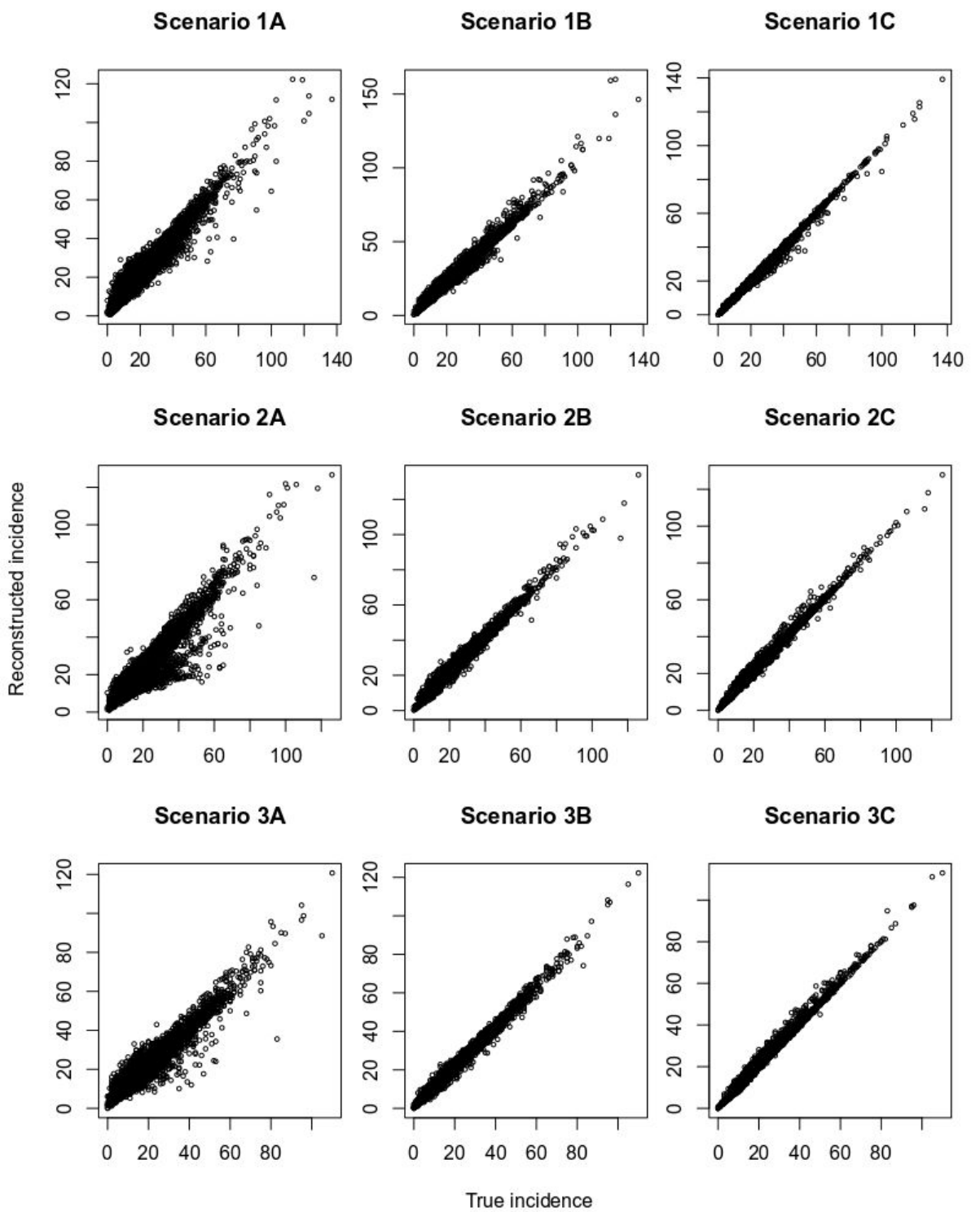
Scenario I _{2A}				Scenario I _{2B}			Scenario I _{2C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.92 (2.87, 2.96)	0.01	2.90	2.94 (2.90, 2.97)	0.01	2.90	2.92 (2.87, 2.97)	0.01
HUM	0.20	0.23 (0.21, 0.25)	0.16	0.20	0.16 (0.14, 0.18)	-0.22	0.20	0.16 (0.14, 0.17)	-0.22
NDVI	0.30	0.27 (0.25, 0.28)	-0.11	0.30	0.30 (0.28, 0.31)	0.00	0.30	0.29 (0.27, 0.30)	-0.04
RAIN	0.20	0.16 (0.14, 0.18)	-0.19	0.20	0.18 (0.16, 0.20)	-0.09	0.20	0.15 (0.13, 0.17)	-0.23
TEMP	0.25	0.18 (0.16, 0.20)	-0.30	0.25	0.22 (0.20, 0.24)	-0.13	0.25	0.21 (0.19, 0.23)	-0.15
IR	0.50	0.49 (0.47, 0.51)	-0.02	0.50	0.46 (0.42, 0.48)	-0.08	0.50	0.47 (0.45, 0.49)	-0.06
σ	10.00	9.47 (9.23, 10.01)	-0.04	15.00	14.85 (14.79, 14.9)	-0.01	20.00	20.26 (19.98, 20.34)	0.01

Scenario I _{3A}				Scenario I _{3B}			Scenario I _{3C}		
Variable	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB	θ	$\hat{\theta}$ (95% CI)	RB
Intercept	2.90	2.94 (2.89, 2.99)	0.01	2.90	2.94 (2.88, 2.98)	0.01	2.90	2.92 (2.87, 2.97)	0.01
HUM	0.20	0.2 (0.18, 0.22)	-0.01	0.20	0.19 (0.18, 0.21)	-0.03	0.20	0.20 (0.18, 0.21)	-0.02
NDVI	0.30	0.3 (0.28, 0.31)	0.00	0.30	0.30 (0.28, 0.31)	-0.01	0.30	0.30 (0.28, 0.31)	-0.01
RAIN	0.20	0.18 (0.16, 0.20)	-0.10	0.20	0.18 (0.16, 0.20)	-0.12	0.20	0.17 (0.15, 0.19)	-0.15
TEMP	0.25	0.20 (0.18, 0.22)	-0.18	0.25	0.23 (0.21, 0.24)	-0.09	0.25	0.23 (0.21, 0.25)	-0.10
IR	0.50	0.46 (0.43, 0.48)	-0.08	0.50	0.48 (0.45, 0.49)	-0.05	0.50	0.48 (0.46, 0.50)	-0.03
σ	10.00	9.76 (9.72, 10.81)	-0.02	15.00	14.48 (14.42, 15.00)	-0.03	20.00	19.55 (19.27, 20.03)	-0.02

Achieving depth and breadth in spatial models of vector-borne diseases: an integrated framework for active survey and passive surveillance data.

Figure 1 displays nine scatter plots arranged in a 3x3 grid, showing the relationship between Reconstructed incidence (Y-axis) and True incidence (X-axis) for various scenarios. The scenarios are labeled as Scenario 1A, Scenario 1B, Scenario 1C, Scenario 2A, Scenario 2B, Scenario 2C, Scenario 3A, Scenario 3B, and Scenario 3C. Each plot shows a dense cluster of points, indicating a strong positive correlation between the reconstructed and true incidence values. The axes are labeled 'Reconstructed incidence' (Y-axis) and 'True incidence' (X-axis). The data points are represented by small open circles.

S3.2 – Plots as showing the relationship between the simulated and reconstructed malaria incidence by model 3.



S3.3 – Plots as showing the relationship between the simulated and reconstructed insecticide resistance by model 3.

